

Spatial adaptation in heteroscedastic regression: Propagation approach

Nora Serdyukova

May 13, 2011

Abstract

The paper concerns the problem of pointwise adaptive estimation in regression when the noise is heteroscedastic and incorrectly known. We use the method of local approximation including as a particular case the local polynomial smoothing. Specifically, the model with unknown mean and variance is approximated by a local linear model with an incorrectly specified covariance matrix.

Adaptive choice of degree of localization in this case can be understood as a choice of an appropriate parametric model from a given collection. For the selection from the family of models we employ based on Lepski's method the FLL technique recently suggested in Katkovnik and Spokoiny (2008). The problem of the choice of certain parameters in this type of procedures was addressed in Spokoiny and Vial (2009). The authors called their approach to the calibration of the parameters "propagation". We developed and justified the methodology for the heteroscedastic case in the presence of noise misspecification. The analysis shows that the adaptive procedure allows a misspecification of the covariance matrix with a relative error of order $(\log n)^{-1}$, where n is the sample size.

1 Introduction

Consider a regression model

$$\mathbf{Y} = \mathbf{f} + \Sigma_0^{1/2} \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(0, I_n) \quad (1.1)$$

with response vector $\mathbf{Y} \in \mathbb{R}^n$ and an unknown covariance matrix $\Sigma_0 = \text{diag}(\sigma_{0,1}^2, \dots, \sigma_{0,n}^2)$. This model can be written as

$$Y_i = f(X_i) + \sigma_{0,i} \varepsilon_i, \quad i = 1, \dots, n \quad (1.2)$$

with design points $X_i \in \mathcal{X} \subset \mathbb{R}^d$. Given a point $x \in \mathcal{X}$, the target of estimation is the value of the regression function $f(x)$ at the point x .

The idea is to replace model (1.2) by a local parametric model

$$y_i = f_{\boldsymbol{\theta}}(X_i) + \sigma_i \varepsilon_i, \quad i : X_i \in U_h(x), \quad (1.3)$$

where $U_h(x) \stackrel{\text{def}}{=} \{t : \|t - x\| \leq h\}$ and $\boldsymbol{\theta} \in \Theta \subset \mathbb{R}^p$ is an unknown finite-dimensional parameter. Then employing one of the well-developed parametric methods we can estimate $\boldsymbol{\theta}$ by $\tilde{\boldsymbol{\theta}}(y_1, \dots, y_d; x)$, and then use the estimator $\tilde{f}_{\tilde{\boldsymbol{\theta}}(Y_1, \dots, Y_d)}(x)$ based on the observations from the

“true” model (1.2) for estimation of $f(x)$. Therefore we have to choose the local model (the collection of estimators $\{f_{\boldsymbol{\theta}}(\cdot), \boldsymbol{\theta} \in \Theta\}$) and the appropriate degree of locality h . This method of local approximation originated from [14], [15], [16], [17]. In what follows we shall consider approximation by local linear models of the following type:

$$y_i = \Psi_i \boldsymbol{\theta} + \sigma_i \epsilon_i, \quad i : X_i \in U_h(x), \quad (1.4)$$

where $\Psi_i = \Psi(X_i) = (\psi_1(X_i - x), \dots, \psi_p(X_i - x))^{\top}$ is a vector of basis functions $\{\psi_j(\cdot)\}$ which already are fixed. The main issue then is the choice of the appropriate bandwidth h such that the estimator

$$\tilde{f}_{\tilde{\boldsymbol{\theta}}_h}(x) \stackrel{\text{def}}{=} \sum_{j=1}^p \tilde{\theta}_h^{(j)}(x) \psi_j(0) \quad (1.5)$$

built on the base of the localized data would be a relevant estimator for $f(x)$. For this purposes the bandwidths selection should be done in a data-driven way, and the adaptive selection from the family $\{\tilde{f}_{\tilde{\boldsymbol{\theta}}_h}(\cdot)\}_{h>0}$ for fixed basis is equivalent to the adaptive choice of bandwidth. Notice also that the coefficients $\theta^{(j)}(x)$ as well as their estimators depends on x and should be calculated for every particular point of interest x . On the other side the localization reduces influence of the choice of the functions $\{\psi_j(\cdot)\}$ allowing to use simple collections.

Moreover, in our set-up the covariance matrix Σ_0 is not assumed to be known exactly, and the approximate model used instead of the true one reads as follows:

$$\mathbf{y} = \boldsymbol{\Psi}^{\top} \boldsymbol{\theta} + \Sigma^{1/2} \boldsymbol{\varepsilon}, \quad (1.6)$$

where $\boldsymbol{\Psi} = (\Psi_1, \dots, \Psi_n)$ is a $p \times n$ design matrix and $\Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_n^2)$, $\min\{\sigma_i^2\} > 0$ is available to a statistician. Thus the model is misspecified in two places: in the form of the regression function and in the error distribution.

The proposed approach includes the important class of polynomial regressions, see [6], [17], [24], [28]. For example in the univariate case $x \in \mathbb{R}$, due to the Taylor theorem, the approximation of the unknown function $f(t)$ for t close to x can be written in the following form: $f_{\boldsymbol{\theta}}(t) = \theta^{(0)} + \theta^{(1)}(t - x) + \dots + \theta^{(p-1)}(t - x)^{p-1}/(p-1)!$ with the parameter $\boldsymbol{\theta} = (\theta^{(0)}, \theta^{(1)}, \dots, \theta^{(p-1)})^{\top}$ corresponding to the values of $f(\cdot)$ and its derivatives at the point x , if they exist. The matrix $\boldsymbol{\Psi}$ then consists of the columns $\Psi_i = (1, X_i - x, \dots, (X_i - x)^{p-1}/(p-1)!)^{\top}$ and corresponds to the well known polynomial smoothing. If the regression function is sufficiently smooth, then for any t close to x , up to a reminder term $f(t) \approx f_{\boldsymbol{\theta}}(t)$, and the estimate of $f(x)$ at the point x is given by $\tilde{f}(x) = \tilde{f}_{\tilde{\boldsymbol{\theta}}(x)}(x) = \tilde{\theta}^{(0)}$. See for more details on local polynomial estimatin [6] or [28]. The local constant fit at a given point $x \in \mathbb{R}$ is covered as well. In this case the design matrix $\boldsymbol{\Psi} = (1, \dots, 1)$, and $f_{\boldsymbol{\theta}}(X_i) = \boldsymbol{\Psi}_i^{\top} \boldsymbol{\theta} = \theta^{(0)} = f_{\boldsymbol{\theta}}(x)$, $i = 1, \dots, n$. This type of approximation in our set-up with known constant noise is treated in [18] and [27].

Nonparametric estimation in heteroscedastic regression under the L_2 losses was studied in [12], [13] and series of papers [8], [9], [10]. For estimation of the mean with L_2 -risk in Gaussian homoscedastic model with unknown variance the penalties allowing to deal with the complexity of such a collection of models were proposed in [2]. However the problem of “local model selection” addressed in the present paper is quite different to the model selection in the sense of [3] and [25] related to estimation with global risk. The minimax pointwise estimation in heteroscedastic regression is in focus of [4].

2 Estimation procedure

2.1 Local parametric estimation

We shall perform the adaptive selection from a collection of K estimators corresponding to model (1.4) with different sizes h of neighborhood $U_h(x)$. Fix a point $x \in \mathbb{R}^d$ as a center of localization and a basis $\{\psi_j\}$. Let the localizing operator be identified with the corresponding matrix. For the next nonparametric step we will need a sequence of nested windows. Thus for every x the sequence of localizing schemes (scales) $\mathcal{W}_k(x)$, $k = 1, \dots, K$ is given by the matrices $\mathcal{W}_k(x) = \text{diag}(w_{k,1}(x), \dots, w_{k,n}(x))$, where the weights $w_{k,i}(x) \in [0, 1]$ can be understood, for instance, as smoothing kernels $w_{k,i}(x) = W((X_i - x)h_k^{-1})$. Let a particular localizing function $w_{(\cdot,\cdot)}(x)$ be fixed; the aim is to choose on the base of available data the index k of the optimal bandwidth h_k . To simplify the notation we sometimes suppress the dependence on the reference point x . Denote by

$$\mathbf{W}_k \stackrel{\text{def}}{=} \Sigma^{-1/2} \mathcal{W}_k \Sigma^{-1/2} = \text{diag}\left(\frac{w_{k,1}}{\sigma_1^2}, \dots, \frac{w_{k,n}}{\sigma_n^2}\right), \quad k = 1, \dots, K. \quad (2.1)$$

Let Θ be a compact subset of \mathbb{R}^p . Inside of any “window” given by \mathbf{W}_k , $k = 1, \dots, K$ we calculate the quasi-maximum likelihood estimator (QMLE) $\tilde{\boldsymbol{\theta}}_k = \tilde{\boldsymbol{\theta}}_k(x) = (\tilde{\theta}_k^{(0)}(x), \dots, \tilde{\theta}_k^{(p-1)}(x))^{\top}$ of $\boldsymbol{\theta}$ defined as

$$\tilde{\boldsymbol{\theta}}_k \stackrel{\text{def}}{=} \underset{\boldsymbol{\theta} \in \Theta}{\text{argmax}} L(\mathbf{W}_k, \boldsymbol{\theta}), \quad (2.2)$$

where

$$\begin{aligned} L(\mathbf{W}_k, \boldsymbol{\theta}) &= -\frac{1}{2} \left(\mathbf{Y} - \boldsymbol{\Psi}^{\top} \boldsymbol{\theta} \right)^{\top} \mathbf{W}_k \left(\mathbf{Y} - \boldsymbol{\Psi}^{\top} \boldsymbol{\theta} \right) + R \\ &= -\frac{1}{2} \sum_{i=1}^n |Y_i - \boldsymbol{\Psi}_i^{\top} \boldsymbol{\theta}|^2 \frac{w_{k,i}}{\sigma_i^2} + R; \end{aligned} \quad (2.3)$$

R stands for the terms not depending on $\boldsymbol{\theta}$, and

$$\boldsymbol{\Psi}_i = \boldsymbol{\Psi}(X_i - x) = (\psi_1(X_i - x), \dots, \psi_p(X_i - x))^{\top}.$$

If the $p \times p$ matrix $\mathbf{B}_k = \mathbf{B}_k(x)$ given by

$$\mathbf{B}_k \stackrel{\text{def}}{=} \boldsymbol{\Psi} \mathbf{W}_k \boldsymbol{\Psi}^{\top} = \sum_{i=1}^n \boldsymbol{\Psi}_i \boldsymbol{\Psi}_i^{\top} \frac{w_{k,i}}{\sigma_i^2} \quad (2.4)$$

is positive definite ($\mathbf{B}_k \succ 0$), then

$$\tilde{\boldsymbol{\theta}}_k = \mathbf{B}_k^{-1} \boldsymbol{\Psi} \mathbf{W}_k \mathbf{Y} = \mathbf{B}_k^{-1} \sum_{i=1}^n \boldsymbol{\Psi}_i Y_i \frac{w_{k,i}}{\sigma_i^2}. \quad (2.5)$$

Recall that in the case of the polynomial basis the estimator $\tilde{\boldsymbol{\theta}}_k(x)$ is a local polynomial estimator of $\boldsymbol{\theta}(x)$ corresponding to the k th scale. In what follows we assume that $n > p$, and that $\det \mathbf{B}_k > 0$ for any $k = 1, \dots, K$. Because $p = \text{rank}(\mathbf{B}_k) \leq \min\{p, \text{rank}(\mathcal{W}_k(x))\}$ this requires the following conditions on the design matrix $\boldsymbol{\Psi}$ and the minimal localizing scheme $\mathcal{W}_1(x)$:

(A1) The $p \times n$ design matrix Ψ is supposed to have full row rank, i.e.,

$$\dim \mathcal{C}(\Psi^\top) = \dim \mathcal{C}(\Psi^\top \Psi) = p.$$

(A2) The smallest localizing scheme $\mathcal{W}_1(x)$ is chosen to contain at least p design points such that $w_{1,i}(x) > 0$, i.e., $p \leq \#\{i : w_{1,i}(x) > 0\}$.

Assumption (A2) in practise is automatically fulfilled, since, for example, in \mathbb{R}^1 it means that for the local constant fit we need at least one observation and so on. Usually it is intrinsically assumed that, starting from the smallest window, at every step of the procedure every new window contains at least p new design points.

The formulas (2.2) give a sequence of estimators $\{\tilde{\theta}_k(x)\}_{k=1}^K$. It was noticed in [1] that in the case when the true data distribution is unknown the QMLE is a natural estimator for the parameter maximizing the expected log-likelihood. That is, for every $k = 1, \dots, K$, the estimator $\tilde{\theta}_k(x)$ can be considered as an estimator of

$$\theta_k^*(x) \stackrel{\text{def}}{=} \underset{\theta \in \Theta}{\operatorname{argmax}} \mathbb{E} L(\mathbf{W}_k, \theta) \quad (2.6)$$

$$\begin{aligned} &= \underset{\theta \in \Theta}{\operatorname{argmin}} (\mathbf{f} - \Psi^\top \theta)^\top \mathbf{W}_k (\mathbf{f} - \Psi^\top \theta) \\ &= \mathbf{B}_k^{-1} \Psi \mathbf{W}_k \mathbf{f} = \mathbf{B}_k^{-1} \sum_{i=1}^n \Psi_i f(X_i) \frac{w_{k,i}}{\sigma_i^2}. \end{aligned} \quad (2.7)$$

Recall that we do not assume $\mathbf{f} = \Psi^\top \theta$ even locally. It is known from [29] that in the presence of a model misspecification for every k the QMLE $\tilde{\theta}_k$ is a strongly consistent estimator for $\theta_k^*(x)$, which also is the minimizer of the localized Kullback-Leibler [19] information criterion:

$$\begin{aligned} \theta_k^*(x) &= \underset{\theta \in \Theta}{\operatorname{argmin}} \sum_{i=1}^n \mathbb{KL}(\mathcal{N}(f(X_i), \sigma_i), \mathcal{N}(\Psi_i^\top \theta, \sigma_i)) w_{k,i}(x) \\ &= \underset{\theta \in \Theta}{\operatorname{argmin}} \sum_{i=1}^n |f(X_i) - \Psi_i^\top \theta|^2 \frac{w_{k,i}(x)}{\sigma_i^2} \end{aligned}$$

with $\mathbb{KL}(P, P_\theta) \stackrel{\text{def}}{=} \mathbb{E}_P [\log(dP/dP_\theta)]$. For the properties of the Kullback-Leibler divergence see, for example, [28].

It follows from the above definition of $\theta_k^*(x)$ and from (2.2) that the QMLE $\tilde{\theta}_k$ admits a decomposition into deterministic and stochastic parts:

$$\tilde{\theta}_k = \mathbf{B}_k^{-1} \Psi \mathbf{W}_k (\mathbf{f} + \Sigma_0^{1/2} \varepsilon) = \theta_k^* + \mathbf{B}_k^{-1} \Psi \mathbf{W}_k \Sigma_0^{1/2} \varepsilon \quad (2.8)$$

$$\mathbb{E} \tilde{\theta}_k = \theta_k^*, \quad (2.9)$$

where $\varepsilon \sim \mathcal{N}(0, I_n)$. Notice that if $\mathbf{f} \equiv \Psi^\top \theta$, then $\theta_k^* \equiv \theta$ for any k , and the classical parametric set-up takes place.

2.2 Adaptive bandwidth selection

Let a point $x \in \mathcal{X} \subset \mathbb{R}^n$, a basis $\{\psi_j\}$ and the method of localization $w_{(.,.)}(x)$ be fixed. The crucial assumption for the procedure under consideration to work is that the localizing schemes

(scales) $\mathcal{W}_k(x) = \text{diag}(w_{k,1}, \dots, w_{k,n})$ are nested. One can say that the localizing schemes are nested in the sense that for the corresponding matrices the following *ordering condition* is fulfilled:

(A3) *For any fixed x and the method of localization $w_{(\cdot,\cdot)}(x)$ the following relation holds:*

$$\mathcal{W}_1(x) \leq \dots \leq \mathcal{W}_k(x) \leq \dots \leq \mathcal{W}_K(x).$$

For the kernel smoothing this condition means the following. Let the sequence of bandwidths $\{h_k\}$ be ordered from the smallest to the largest one, i.e., $h_1 < \dots < h_K$, and let $\mathcal{W}_k(x) = \text{diag}(w_{k,1}, \dots, w_{k,n})$ be the localizing matrix, corresponding to the bandwidth h_k . Here the weights $w_{k,i} = w_{k,i}(x) = W((X_i - x)/h_k) \in [0, 1]$ are nonnegative functions such that $W(u/h_l) \leq W(u/h_k)$ for any $0 < h_l < h_k < 1$, and $W(u) \rightarrow 0$ as $|u| \rightarrow \infty$, or even is compactly supported.

Recall that, given $x \in \mathcal{X}$, a basis $\{\psi_j\}$, and the method of localization $w_{(\cdot,\cdot)}(x)$, we look for the estimator $\hat{f}_{\hat{\theta}}(x)$ of $f(x)$ having form (1.5), where the coefficients $\hat{\theta}^{(j)}(x)$ are the components of the estimator

$$\hat{\theta}(x) \stackrel{\text{def}}{=} \tilde{\theta}_{\hat{k}}(x) = (\tilde{\theta}_{\hat{k}}^{(1)}(x), \dots, \tilde{\theta}_{\hat{k}}^{(p)}(x))^{\top}, \quad (2.10)$$

corresponding to the adaptive choice of the index $\hat{k} \in \{1, \dots, K\}$, i.e. to the choice of the degree of localization.

The selection of $\hat{\theta}(x)$ from $\{\tilde{\theta}_k(x)\}$, $k = 1, \dots, K$ can be done by the application of the Lepski [20] method to the comparing of the maximized log-likelihoods $L(\mathbf{W}_k, \tilde{\theta}_k)$. This is the idea of the FLL technique suggested in [18]. More precisely, to describe the test statistic, define for any $\theta, \theta' \in \Theta$ the corresponding log-likelihood ratio:

$$L(\mathbf{W}_k, \theta, \theta') \stackrel{\text{def}}{=} L(\mathbf{W}_k, \theta) - L(\mathbf{W}_k, \theta'). \quad (2.11)$$

Then, using the approach suggested in [18], for every $l = 1, \dots, K$, the *fitted log-likelihood (FLL)* ratio is defined as follows:

$$L(\mathbf{W}_l, \tilde{\theta}_l, \theta') \stackrel{\text{def}}{=} \max_{\theta \in \Theta} L(\mathbf{W}_l, \theta, \theta').$$

By Theorem 4.2, for any l and θ , the FLL is a quadratic form:

$$2L(\mathbf{W}_l, \tilde{\theta}_l, \theta) = (\tilde{\theta}_l - \theta)^{\top} \mathbf{B}_l (\tilde{\theta}_l - \theta).$$

This prompts the use, see [18], the *FLL-statistics*:

$$\begin{aligned} T_{lk} &\stackrel{\text{def}}{=} 2L(\mathbf{W}_l, \tilde{\theta}_l, \tilde{\theta}_k) \\ &= (\tilde{\theta}_l - \tilde{\theta}_k)^{\top} \mathbf{B}_l (\tilde{\theta}_l - \tilde{\theta}_k), \quad l < k. \end{aligned} \quad (2.12)$$

In the algorithm the smallest bandwidths corresponding to $k = 1$ is always accepted, then the adaptive index \hat{k} is selected by Lepski's selection rule with the FLL test statistics $\{T_{lm}\}$, $1 \leq l < m \leq K$:

$$\hat{k} = \max \{k \leq K : T_{lm} \leq \zeta_l, l < m \leq k\}. \quad (2.13)$$

Finally we set $\hat{\theta} = \tilde{\theta}_{\hat{k}}$.

The procedure (2.13) involves parameters $\zeta_1, \dots, \zeta_{K-1}$ related to the large deviations of $\{T_{lm}\}$, $1 \leq l < m \leq K$. As the classical Lepski procedure, the (2.13) controls the risk of estimators for the case of dominating bias. The opposite case of the negligible w.r.t. the noise bias is usually

handled by employing the advanced empirical process technique, however sometimes providing the constants far away of being optimal. Notice also that the Wilks-type Theorem 4.3 below gives the bound for the expected fitted log-likelihood ratio:

$$\mathbb{E}|2L(\mathbf{W}_k, \tilde{\boldsymbol{\theta}}_k, \boldsymbol{\theta}_k^*)|^r \leq C(p, r) \quad (2.14)$$

where the constant $C(p, r)$ does not depend on the degree of localization and is given by the r th moment of the χ^2 distribution with p degrees of freedom:

$$C(p, r) = \mathbb{E}|\chi_p^2|^r = 2^r \Gamma(r + p/2)/\Gamma(p/2), \quad (2.15)$$

Therefore we shall follow the practical idea from [27] and [18] allowing to avoid hard large deviations analysis and to calculate the thresholds rather sharp numerically. We assume at this step that the critical values $\mathfrak{z}_1, \dots, \mathfrak{z}_{K-1}$ are already fixed satisfying the following set of $K - 1$ inequalities:

Definition 2.1. *Propagation conditions (PC)*

Let $\hat{\boldsymbol{\theta}}_k$ denote the last accepted estimate after the first k steps of the procedure:

$$\hat{\boldsymbol{\theta}}_k \stackrel{\text{def}}{=} \tilde{\boldsymbol{\theta}}_{\min\{k, \hat{k}\}}. \quad (2.16)$$

The critical values $\mathfrak{z}_1, \dots, \mathfrak{z}_{K-1}$ satisfy

$$\mathbb{E}_{0, \Sigma} |(\tilde{\boldsymbol{\theta}}_k - \hat{\boldsymbol{\theta}}_k)^\top \mathbf{B}_k (\tilde{\boldsymbol{\theta}}_k - \hat{\boldsymbol{\theta}}_k)|^r \leq \alpha C(p, r) \quad \text{for all } k = 2, \dots, K, \quad (2.17)$$

where $C(p, r)$ is defined by (2.15), $\alpha \in (0, 1]$ is an additional “free” turning parameter which can be taken equal to 1, and $\mathbb{E}_{0, \Sigma}$ stands for the expectation w.r.t. the measure $\mathcal{N}(0, \Sigma)$.

Remark 2.1. Lemma 4.1 from Section 4 shows that in the “no bias” situation the Gaussian distribution provides a nice pivotality property: the actual value of the parameter $\boldsymbol{\theta}$ is not important for the risk of adaptive estimate, so one can put $\boldsymbol{\theta} = 0$ in (2.17).

Remark 2.2. Clearly at any step $k \leq K$ of the algorithm the “current value” of the adaptive estimator $\hat{\boldsymbol{\theta}}_k$ depends on the thresholds $\mathfrak{z}_1, \dots, \mathfrak{z}_{k-1}$. The theoretical aspects related to the heteroscedasticity of model and to the incorrectly known variance is the focus of the present paper. Thus we do not detail the practical aspects of the thresholds calibration only mentioning that in practise this can be done by Monte Carlo simulations under the known “parametric” measure $\mathcal{N}(0, \Sigma)$. Moreover one needs to calculate them only once. For detailed consideration of the practical aspects of the calibration as well as for the computational results see [27] or [18] focused on the image denoising by local constant fitting. Demo-versions of the software are available on the web page <http://www.cs.tut.fi/~lasip/>.

3 Theoretical study

In order to control the admissible level of misspecification for “model” covariance matrix from (1.6) we need to introduce the following condition on the relative variability in errors:

(A4) There exists $\delta \in [0, 1)$ such that

$$1 - \delta \leq \sigma_{0,i}^2 / \sigma_i^2 \leq 1 + \delta \quad \text{for all } i = 1, \dots, n.$$

3.1 Upper bound for the critical values

Let us at this step recall the notion of the *Löwner partial ordering*: for any real symmetric matrices A and B we write $A \preceq B$ if and only if $\vartheta^\top A \vartheta \leq \vartheta^\top B \vartheta$ for all vectors ϑ , or, equivalently, if and only if the matrix $B - A$ is nonnegative definite. Assuming (A4), the true covariance matrix $\Sigma_0 \preceq \Sigma(1 + \delta)$, and the variance of the estimate $\tilde{\boldsymbol{\theta}}_k$ is bounded with \mathbf{B}_k^{-1} :

$$V_k \stackrel{\text{def}}{=} \text{Var} \tilde{\boldsymbol{\theta}}_k = \mathbf{B}_k^{-1} \Psi \mathbf{W}_k \Sigma_0 \mathbf{W}_k \Psi^\top \mathbf{B}_k^{-1} \quad (3.1)$$

$$\begin{aligned} &\preceq (1 + \delta) \mathbf{B}_k^{-1} \Psi \mathbf{W}_k \Sigma \mathbf{W}_k \Psi^\top \mathbf{B}_k^{-1} \\ &= (1 + \delta) \mathbf{B}_k^{-1} \Psi \Sigma^{-1/2} \mathcal{W}_k^2 \Sigma^{-1/2} \Psi^\top \mathbf{B}_k^{-1} \\ &\preceq (1 + \delta) \mathbf{B}_k^{-1} \Psi \Sigma^{-1/2} \mathcal{W}_k \Sigma^{-1/2} \Psi^\top \mathbf{B}_k^{-1} \\ &= (1 + \delta) \mathbf{B}_k^{-1} \Psi \mathbf{W}_k \Psi^\top \mathbf{B}_k^{-1} \\ &= (1 + \delta) \mathbf{B}_k^{-1}. \end{aligned} \quad (3.2)$$

The last inequality follows from the observation that all the entries of the “weight” matrix \mathcal{W}_k do not exceed one, implying $\mathcal{W}_k^2 \preceq \mathcal{W}_k$. The strict equality takes place if the $\{w_{k,i}\}$ are boxcar (rectangular) kernels and the noise is known, i.e., if $\delta = 0$. To justify the procedure one need to show that the critical values chosen by (PC) are finite. This is obtained under the following assumption:

(A5) *Let for some constants u_0 and u such that $1 < u_0 \leq u$ for any $2 \leq k \leq K$ the matrices \mathbf{B}_k satisfy*

$$u_0 I_p \preceq \mathbf{B}_{k-1}^{-1/2} \mathbf{B}_k \mathbf{B}_{k-1}^{-1/2} \preceq u I_p$$

Remark 3.1. In the “one dimensional case” $p = 1$, that is for the local constant fit, the “matrix” $\mathbf{B}_k = \sum_{i=1}^n w_{k,i} \sigma_i^{-2} \geq \mathbf{B}_{k-1}$ is just a weighted “local design size”. Assume for simplicity that $\sigma_i^2 \equiv \sigma^2$, the weights are rectangular kernels $w_{k,i}(x) = \mathbb{I}\{|X_i - x| \leq h_k/2\}$, and the design is equidistant. Then for n sufficiently large

$$\frac{1}{n} \mathbf{B}_k = \frac{1}{n\sigma^2} \sum_{i=1}^n \mathbb{I}\left\{ \left| \frac{i}{n} - x \right| \leq \frac{h_k}{2} \right\} \approx \frac{h_k}{\sigma^2},$$

and the condition (A5) means that the bandwidths grow geometrically: $h_k = u h_{k-1}$.

Now we are able to formulate a theorem on finiteness of the critical values.

Theorem 3.1. The theoretical choice of the critical values

Assume (A1) – (A3) and (A5). The adaptive procedure (2.13) in the considered set-up is well defined in the sense that the choice of the critical values of the form

$$\mathfrak{z}_k = \frac{4}{\mu} \left\{ r(K - k) \log u + \log(K/\alpha) - \frac{p}{4} \log(1 - 4\mu) - \log(1 - u^{-r}) + \overline{C}(p, r) \right\} \quad (3.3)$$

provides the conditions (2.17) for all $k \leq K$. Here

$$\overline{C}(p, r) = \log \left\{ \frac{2^{2r} [\Gamma(2r + p/2)\Gamma(p/2)]^{1/2}}{\Gamma(r + p/2)} \right\}$$

and $\mu \in (0, 1/4)$. Particularly,

$$\mathbb{E}_{0,\Sigma} |(\tilde{\boldsymbol{\theta}}_K - \hat{\boldsymbol{\theta}})^\top \mathbf{B}_K (\tilde{\boldsymbol{\theta}}_K - \hat{\boldsymbol{\theta}})|^r \leq \alpha C(p, r). \quad (3.4)$$

The proof is given in subsection 4.2.

3.2 Quality of estimation in the nearly parametric case: The small modeling bias condition and the propagation property

The critical values of the procedure $\mathfrak{z}_1, \dots, \mathfrak{z}_{K-1}$ were selected by the propagation conditions (2.17) under the measure $\mathcal{N}(\boldsymbol{\theta}, \Sigma)$. Now $\boldsymbol{\theta}_1^* \approx \dots \approx \boldsymbol{\theta}_k^* \approx \boldsymbol{\theta}$ up to some $k \leq K$, and the covariance matrix is Σ_0 . The aim is to formalize the meaning of “ \approx ” and to justify the use of the critical values in this situation. For this purposes we will take into account the discrepancy between the joint distributions of linear estimates $\tilde{\boldsymbol{\theta}}_1, \dots, \tilde{\boldsymbol{\theta}}_k$ for $k = 1, \dots, K$ under “no bias” assumption corresponding to the distributions with the mean $\boldsymbol{\theta}_1^* = \dots = \boldsymbol{\theta}_k^* = \boldsymbol{\theta}$ and the incorrectly specified covariance matrix Σ , and in the general situation with $\boldsymbol{\theta}_1^* \neq \dots \neq \boldsymbol{\theta}_k^*$ and the covariance Σ_0 . Denote the expectations w.r.t. these measures by $\mathbb{E}_{\boldsymbol{\theta}, \Sigma} := \mathbb{E}_{k, \boldsymbol{\theta}, \Sigma}$ and $\mathbb{E}_{\boldsymbol{f}, \Sigma_0} := \mathbb{E}_{k, \boldsymbol{f}, \Sigma_0}$, respectively. Denote a $p \times k$ matrix of the first k estimators and the expectations correspondingly by

$$\begin{aligned}\tilde{\boldsymbol{\Theta}}_k &\stackrel{\text{def}}{=} (\tilde{\boldsymbol{\theta}}_1, \dots, \tilde{\boldsymbol{\theta}}_k), \\ \boldsymbol{\Theta}_k^* &\stackrel{\text{def}}{=} \mathbb{E}_{\boldsymbol{f}, \Sigma_0} \tilde{\boldsymbol{\Theta}}_k = (\boldsymbol{\theta}_1^*, \dots, \boldsymbol{\theta}_k^*), \\ \boldsymbol{\Theta}_k &\stackrel{\text{def}}{=} \mathbb{E}_{\boldsymbol{\theta}, \Sigma} \tilde{\boldsymbol{\Theta}}_k = (\boldsymbol{\theta}, \dots, \boldsymbol{\theta}).\end{aligned}$$

Let $A \otimes B$ stands for the Kronecker product of A and B defined as

$$A \otimes B = \begin{pmatrix} a_{11}B & a_{12}B & \cdots & a_{1n}B \\ a_{21}B & a_{22}B & \cdots & a_{2n}B \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1}B & a_{m2}B & \cdots & a_{mn}B \end{pmatrix}.$$

Denote the $pk \times pk$ covariance matrices of $\text{vec } \tilde{\boldsymbol{\Theta}}_k^\top = (\tilde{\boldsymbol{\theta}}_1^\top, \dots, \tilde{\boldsymbol{\theta}}_k^\top) \in \mathbb{R}^{pk}$ by

$$\boldsymbol{\Sigma}_k \stackrel{\text{def}}{=} \text{Var}_{\boldsymbol{\theta}, \Sigma} [\text{vec } \tilde{\boldsymbol{\Theta}}_k] = \mathbf{D}_k (J_k \otimes \Sigma) \mathbf{D}_k^\top, \quad (3.5)$$

$$\boldsymbol{\Sigma}_{k,0} \stackrel{\text{def}}{=} \text{Var}_{\boldsymbol{f}, \Sigma_0} [\text{vec } \tilde{\boldsymbol{\Theta}}_k] = \mathbf{D}_k (J_k \otimes \Sigma_0) \mathbf{D}_k^\top, \quad (3.6)$$

where the matrix J_k is a $k \times k$ matrix with all its elements equal to 1, and the $pk \times nk$ matrix \mathbf{D}_k is defined as follows:

$$\begin{aligned}\mathbf{D}_k &\stackrel{\text{def}}{=} D_1 \oplus \dots \oplus D_k = \text{diag}(D_1, \dots, D_k), \\ D_l &\stackrel{\text{def}}{=} \mathbf{B}_l^{-1} \boldsymbol{\Psi} \mathbf{W}_l, \quad l = 1, \dots, k.\end{aligned} \quad (3.7)$$

By Lemma 4.8 from Section 4 under Assumption (A4) with the same δ the similar relation holds for the covariance matrices $\boldsymbol{\Sigma}_k$ and $\boldsymbol{\Sigma}_{k,0}$ of the linear estimators:

$$(1 - \delta) \boldsymbol{\Sigma}_k \preceq \boldsymbol{\Sigma}_{k,0} \preceq (1 + \delta) \boldsymbol{\Sigma}_k, \quad k \leq K. \quad (3.8)$$

In spite of the moment generating function of $\text{vec } \tilde{\Theta}_K$ has the form corresponding to the multivariate normal distribution, see Lemma 4.10 in Section 4, this representation makes sense only if Σ_K is nonsingular. Notice that $\text{rank}(J_K \otimes \Sigma) = n$. From $J_K \otimes \Sigma \succeq 0$ it follows only that $\Sigma_K \succeq 0$, similarly, $\Sigma_{K,0} \succeq 0$. However, without any additional assumptions it is easy to show, see Lemma 4.9 in Section 4, that for rectangular kernels $\Sigma_K \succ 0$. On the other hand, due to (3.8), it is enough to require nonsingularity only for the matrix Σ_K corresponding to the approximate model (1.6), and its choice belongs to a statistician. In what follows we assume that $\Sigma_K \succ 0$.

Denote by $\mathbb{P}_{\theta,\Sigma}^k = \mathcal{N}(\text{vec } \Theta_k, \Sigma_k)$ and by $\mathbb{P}_{f,\Sigma_0}^k = \mathcal{N}(\text{vec } \Theta_k^*, \Sigma_{k,0})$, $k = 1, \dots, K$, the distributions of $\text{vec } \tilde{\Theta}_k$ under the null and under the alternative. Denote also the Radon-Nikodym derivative by

$$Z_k \stackrel{\text{def}}{=} \frac{d\mathbb{P}_{f,\Sigma_0}^k}{d\mathbb{P}_{\theta,\Sigma}^k}. \quad (3.9)$$

Then, by Lemma 4.11 from Section 4, the Kullback-Leibler divergence between these measures has the following form:

$$\begin{aligned} 2\mathbb{KL}(\mathbb{P}_{f,\Sigma_0}^k, \mathbb{P}_{\theta,\Sigma}^k) &\stackrel{\text{def}}{=} 2\mathbb{E}_{f,\Sigma_0} \log \left(\frac{d\mathbb{P}_{f,\Sigma_0}^k}{d\mathbb{P}_{\theta,\Sigma}^k} \right) \\ &= \Delta(k) + \log \left(\frac{\det \Sigma_k}{\det \Sigma_{k,0}} \right) + \text{tr}(\Sigma_k^{-1} \Sigma_{k,0}) - pk, \end{aligned} \quad (3.10)$$

where

$$b(k) \stackrel{\text{def}}{=} \text{vec } \Theta_k^* - \text{vec } \Theta_k \quad (3.11)$$

$$\Delta(k) \stackrel{\text{def}}{=} b(k)^\top \Sigma_k^{-1} b(k). \quad (3.12)$$

If there would be no any ‘‘noise misspecification’’, i.e., if $\delta \equiv 0$ implying $\Sigma = \Sigma_0$, then $\Delta(k) = b(k)^\top \Sigma_k^{-1} b(k) = 2\mathbb{KL}(\mathbb{P}_{f,\Sigma}^k, \mathbb{P}_{\theta,\Sigma}^k)$. Therefore, this quantity can be used to indicate deviation between the mean values in the true (1.1) and the approximate (1.6) models. Clearly, under (\mathcal{W}) , the quantity $\Delta(k)$ grows with k , so following the terminology suggested in [27], we introduce the *small modeling bias condition*:

(SMB) *Let for some $k \leq K$ and some θ exist a constant $\Delta \geq 0$ such that*

$$\Delta(k) \leq \Delta.$$

Monotonicity of $\Delta(k)$ and Assumption (SMB) immediately imply that

$$\Delta(k') \leq \Delta \text{ for all } k' \leq k.$$

The conditions (3.8) yield $-pk\delta \leq \text{tr}(\Sigma_k^{-1} \Sigma_{k,0}) - pk \leq pk\delta$. Thus (4.19) implies the bound for the Kullback-Leibler divergence in terms of δ :

$$-\frac{pk}{2} \log(1 + \delta) + \frac{\Delta(k)}{2} - \frac{pk\delta}{2} \leq \mathbb{KL}(\mathbb{P}_{f,\Sigma_0}^k, \mathbb{P}_{\theta,\Sigma}^k) \leq -\frac{pk}{2} \log(1 - \delta) + \frac{\Delta(k)}{2} + \frac{pk\delta}{2}. \quad (3.13)$$

Moreover, as $\delta \rightarrow 0+$

$$\Delta(k) - 2pk\delta + o(\delta) \leq 2\mathbb{KL}(\mathbb{P}_{f,\Sigma_0}^k, \mathbb{P}_{\theta,\Sigma}^k) \leq \Delta(k) + 2pk\delta + o(\delta). \quad (3.14)$$

This means that, if for some k Assumption (SMB) is fulfilled and $\delta = O(1/K)$, then the Kullback-Leibler divergence between the measures $\mathbb{P}_{\theta,\Sigma}^k$ and $\mathbb{P}_{f,\Sigma_0}^k$ is bounded by a small constant.

Now one can state the crucial property for obtaining the final oracle result.

Theorem 3.2. Propagation property

Assume (A1) – (A5) and (PC). Then for any $k \leq K$ the following upper bounds hold:

$$\begin{aligned} & \mathbb{E}|(\tilde{\boldsymbol{\theta}}_k - \boldsymbol{\theta})^\top \mathbf{B}_k (\tilde{\boldsymbol{\theta}}_k - \boldsymbol{\theta})|^{r/2} \\ & \leq (\mathbb{E}|\chi_p^2|^r)^{1/2} (1 + \delta)^{pk/4} (1 - \delta)^{-3pk/4} \exp \left\{ \varphi(\delta) \frac{\Delta(k)}{2(1 - \delta)} \right\}, \\ & \mathbb{E}|(\tilde{\boldsymbol{\theta}}_k - \hat{\boldsymbol{\theta}}_k)^\top \mathbf{B}_k (\tilde{\boldsymbol{\theta}}_k - \hat{\boldsymbol{\theta}}_k)|^{r/2} \\ & \leq (\alpha \mathbb{E}|\chi_p^2|^r)^{1/2} (1 + \delta)^{pk/4} (1 - \delta)^{-3pk/4} \exp \left\{ \varphi(\delta) \frac{\Delta(k)}{2(1 - \delta)} \right\}, \end{aligned}$$

where $\varphi(\delta) \stackrel{\text{def}}{=} \begin{cases} 1 & \text{for homogeneous errors,} \\ \frac{2(1+\delta)}{(1-\delta)^2} - 1 & \text{otherwise.} \end{cases}$

Here $\tilde{\boldsymbol{\theta}}_k = \tilde{\boldsymbol{\theta}}_k(x)$ is the QMLE defined by (2.2), and $\hat{\boldsymbol{\theta}}_k(x) = \tilde{\boldsymbol{\theta}}_{\min k, \hat{k}}(x)$ is the adaptive estimate at the k th step of the procedure.

The proof is given in Subsection 4.4.

Remark 3.2. Bounds (4.28) and (4.27) below give a condition on the relative error in the noise misspecification. As $\delta \rightarrow 0+$ for every $k \leq K$ it holds that

$$\varphi(\delta) \frac{\Delta(k)}{1 + \delta} - 2pk\delta + o(\delta) \leq \log \mathbb{E}_{\boldsymbol{\theta}, \Sigma}[Z_k^2] \leq \varphi(\delta) \frac{\Delta(k)}{1 - \delta} + 2pk\delta + o(\delta),$$

where Z_k is defined by (3.9). This bound implies, up to the additive constant $\log(\alpha \mathbb{E}|\chi_p^2|^r)/2$, the same asymptotic behavior for the logarithm of the risk of adaptive estimate at each step of the procedure. Because by (SMB) the quantity $\Delta(k)$ is supposed to be bounded by a small constant, and K is of order $\log n$, then $\mathbb{E}_{\boldsymbol{\theta}, \Sigma}[Z_k^2]$ is small if $\delta = O(1/\log n)$. This means that for the case when Σ is an estimate for Σ_0 only the logarithmic in sample size quality is needed. This observation is of particular importance, since it is known from [26] that over classes of functions with bounded second derivative the rate $n^{-1/2}$ of variance estimation is achievable only for the dimension $d \leq 8$.

Remark 3.3. The propagation property provides the adaptive procedure do not stop with high probability while $\Delta(k)$ is small, i.e., under (SMB), and if the relative error δ in the noise is sufficiently small.

3.3 Quality of estimation in the nonparametric case: the oracle result

Define the *oracle index* as the largest index $k \leq K$ such that (SMB) holds:

$$k^* \stackrel{\text{def}}{=} \max\{k \leq K : \Delta(k) \leq \Delta\}. \quad (3.15)$$

Theorem 3.3. Let $\Delta(1) \leq \Delta$, i.e., the first estimate is always accepted in the testing procedure. Let k^* be the oracle index. Then under the conditions (A1), (A2), (A4), (\mathcal{W}), (A5) the risk between the adaptive estimate and the oracle one is bounded with the following expression:

$$\begin{aligned} & \mathbb{E}|(\tilde{\boldsymbol{\theta}}_{k^*} - \hat{\boldsymbol{\theta}})^\top \mathbf{B}_{k^*}(\tilde{\boldsymbol{\theta}}_{k^*} - \hat{\boldsymbol{\theta}})|^{r/2} \\ & \leq \mathfrak{z}_{k^*}^{r/2} + (\alpha \mathbb{E}|\chi_p^2|^r)^{1/2} (1+\delta)^{pk^*/4} (1-\delta)^{-3pk^*/4} \exp\left\{\varphi(\delta) \frac{\Delta}{2(1-\delta)}\right\}, \end{aligned} \quad (3.16)$$

where $\varphi(\delta)$ is as in Theorem 3.2.

Proof. By the definition of the adaptive estimate $\hat{\boldsymbol{\theta}} = \tilde{\boldsymbol{\theta}}_{\hat{k}}$. Because the events $\{\hat{k} \leq k^*\}$ and $\{\hat{k} > k^*\}$ are disjunct, one can write

$$\begin{aligned} & \mathbb{E}|(\tilde{\boldsymbol{\theta}}_{k^*} - \hat{\boldsymbol{\theta}})^\top \mathbf{B}_{k^*}(\tilde{\boldsymbol{\theta}}_{k^*} - \hat{\boldsymbol{\theta}})|^{r/2} \\ &= \mathbb{E}|(\tilde{\boldsymbol{\theta}}_{k^*} - \tilde{\boldsymbol{\theta}}_{\hat{k}})^\top \mathbf{B}_{k^*}(\tilde{\boldsymbol{\theta}}_{k^*} - \tilde{\boldsymbol{\theta}}_{\hat{k}})|^{r/2} \mathbb{I}\{\hat{k} \leq k^*\} \\ &+ \mathbb{E}|(\tilde{\boldsymbol{\theta}}_{k^*} - \tilde{\boldsymbol{\theta}}_{\hat{k}})^\top \mathbf{B}_{k^*}(\tilde{\boldsymbol{\theta}}_{k^*} - \tilde{\boldsymbol{\theta}}_{\hat{k}})|^{r/2} \mathbb{I}\{\hat{k} > k^*\}. \end{aligned}$$

If $\hat{k} \leq k^*$ then $\tilde{\boldsymbol{\theta}}_{k^*} \stackrel{\text{def}}{=} \tilde{\boldsymbol{\theta}}_{\min\{k^*, \hat{k}\}} = \tilde{\boldsymbol{\theta}}_{\hat{k}}$. Thus, to bound the first summand, it is enough to apply Theorem 3.2 with $k = k^*$.

To bound the second expectation, i.e., to bound fluctuations of the adaptive estimate $\hat{\boldsymbol{\theta}}$ at the steps of the procedure for which the (SMB) condition is not fulfilled anymore, just notice that for $\hat{k} > k^*$ the quadratic form coincides with the test statistics $T_{k^*, \hat{k}}$

$$\begin{aligned} & (\tilde{\boldsymbol{\theta}}_{k^*} - \hat{\boldsymbol{\theta}})^\top \mathbf{B}_{k^*}(\tilde{\boldsymbol{\theta}}_{k^*} - \hat{\boldsymbol{\theta}}) \\ &= (\tilde{\boldsymbol{\theta}}_{k^*} - \tilde{\boldsymbol{\theta}}_{\hat{k}})^\top \mathbf{B}_{k^*}(\tilde{\boldsymbol{\theta}}_{k^*} - \tilde{\boldsymbol{\theta}}_{\hat{k}}) \stackrel{\text{def}}{=} T_{k^*, \hat{k}}. \end{aligned}$$

But the index \hat{k} was accepted, this means that $T_{l, \hat{k}} \leq \mathfrak{z}_l$ for all $l < \hat{k}$ and therefore for $l = k^*$. Thus

$$\mathbb{E}|(\tilde{\boldsymbol{\theta}}_{k^*} - \hat{\boldsymbol{\theta}})^\top \mathbf{B}_{k^*}(\tilde{\boldsymbol{\theta}}_{k^*} - \hat{\boldsymbol{\theta}})|^{r/2} \mathbb{I}\{\hat{k} > k^*\} \leq \mathfrak{z}_{k^*}^{r/2}.$$

□

3.4 Componentwise oracle risk bounds

Theorem 3.3 provides the oracle risk bound for the adaptive estimator $\hat{\boldsymbol{\theta}}(x) = \tilde{\boldsymbol{\theta}}_{\hat{k}}(x)$ of the parameter vector $\boldsymbol{\theta}(x) \in \mathbb{R}^p$ corresponding to the estimator $\hat{f}_{\boldsymbol{\theta}}(x)$ of the type (1.5). It is interesting to have a look at the oracle quality of estimation of the components $\theta^{(1)}, \dots, \theta^{(p)}$ of the vector $\boldsymbol{\theta}$ having in mind that the choice of polynomial basis leads to the direct estimation of the value of regression function and the derivatives by the coordinates of $\hat{\boldsymbol{\theta}}$.

Denote by $LP_k(p-1)$ a local polynomial estimator of order $p-1$ corresponding to the k th degree of localization, and, respectively, by $LP^{ad}(p-1)$ its adaptive counterpart, i.e., $LP^{ad}(p-1) \stackrel{\text{def}}{=} LP_{\hat{k}}(p-1)$. If the basis is polynomial and the regression function $f(\cdot)$ is sufficiently smooth in a neighborhood of x , then $\hat{\boldsymbol{\theta}}(x)$ is the $LP^{ad}(p-1)$ of the vector $(f^{(0)}(x), \dots, f^{(p-1)}(x))^\top$ of the values of the function f and its derivatives at the reference point $x \in \mathbb{R}^d$.

Now we are going to obtain a similar to the previous section oracle result for the components of the vector $\widehat{\boldsymbol{\theta}}(x)$, particularly for $\mathbf{e}_j^\top \widehat{\boldsymbol{\theta}}(x)$, $j = 1, \dots, p$, where $\mathbf{e}_j = (0, \dots, 1, \dots, 0)^\top$ is the j th canonical basis vector in \mathbb{R}^p . As a corollary of this general result in the case of the polynomial basis we get an oracle risk bound for $LP^{ad}(p-1)$ estimators of the function f and its derivatives at the point x .

Let $LP_k(p-1)$ estimator of $f^{(j-1)}(x)$ be given by

$$\begin{aligned}\tilde{f}_k^{(j-1)}(x) &= \mathbf{e}_j^\top \widetilde{\boldsymbol{\theta}}_k(x), \quad j = 1, \dots, p, \\ \tilde{f}_k(x) &= \tilde{f}_k^{(0)}(x) = \mathbf{e}_1^\top \widetilde{\boldsymbol{\theta}}_k(x).\end{aligned}\tag{3.17}$$

Then the adaptive local polynomial estimators are defined as follows:

$$\begin{aligned}\widehat{f}^{(j-1)}(x) &= \mathbf{e}_j^\top \widehat{\boldsymbol{\theta}}(x), \quad j = 1, \dots, p, \\ \widehat{f}(x) &= \mathbf{e}_1^\top \widehat{\boldsymbol{\theta}}(x).\end{aligned}\tag{3.18}$$

Similarly, the adaptive estimators of the function f and its derivatives corresponding to the k th step of the procedure are given by

$$\widehat{f}_k^{(j-1)}(x) \stackrel{\text{def}}{=} \mathbf{e}_j^\top \widehat{\boldsymbol{\theta}}_k(x), \quad j = 1, \dots, p.\tag{3.19}$$

Thus, if the basis is polynomial, the estimator $\widehat{f}(x) \stackrel{\text{def}}{=} \widehat{f}^{(0)}(x)$ is the $LP^{ad}(p-1)$ estimator of the value $f(x)$, and $\widehat{f}^{(j-1)}(x)$ with $j = 2, \dots, p$ are, correspondingly, the $LP^{ad}(p-1)$ estimators of the values of its derivatives. However it should be stressed that the results of Theorems 3.3 and 3.9 hold for any basis satisfying the conditions of the theorems. We shall need the following assumptions:

(A6) *There exist $0 < \sigma_{\min}(k) \leq \sigma_{\max}(k) < \infty$ such that for $i : X_i \in U_{h_k}(x)$, with $U_{h_k}(x)$ given by \mathbf{W}_k the variances of errors from the parametric (known) model (1.6) are locally uniformly bounded:*

$$\sigma_{\min}^2(k) \leq \sigma_i^2 \leq \sigma_{\max}^2(k).$$

(A7) *Let assumption (A6) be satisfied. There exists a number $\Lambda_0 > 0$ such that for any $k = 1, \dots, K$ the smallest eigenvalue $\lambda_p(\mathbf{B}_k) \geq nh_k^d \Lambda_0 \sigma_{\max}^{-2}(k)$ for n sufficiently large.*

Then, because $\mathbf{B}_k \succ 0$, for any $k = 1, \dots, K$ and for any $\gamma \in \mathbb{R}^p$ we have

$$\gamma^\top \mathbf{B}_k^{-1} \gamma \leq \frac{\sigma_{\max}^2(k)}{nh_k^d \Lambda_0} \|\gamma\|^2 \leq \frac{\bar{\sigma}_{\max}^2(k)}{nh_k^d \Lambda_0} \|\gamma\|^2,\tag{3.20}$$

where $\bar{\sigma}_{\max}^2(k) \stackrel{\text{def}}{=} \max_{1 \leq l \leq k} \sigma_{\max}^2(l)$. Thus we have the following lemma:

Lemma 3.4. *Let (A6) and (A7) be satisfied. Then for any $j = 1, \dots, p$ and $k, k' = 1, \dots, K$ the following upper bound holds:*

$$\left(\frac{nh_k^d \Lambda_0}{\bar{\sigma}_{\max}^2(k)} \right)^{1/2} |\mathbf{e}_j^\top \widetilde{\boldsymbol{\theta}}_k - \mathbf{e}_j^\top \widetilde{\boldsymbol{\theta}}_{k'}| \leq \|\mathbf{B}_k^{1/2} (\widetilde{\boldsymbol{\theta}}_k - \widetilde{\boldsymbol{\theta}}_{k'})\|.$$

Proof. By (3.20) taking $\gamma = \mathbf{B}_k^{1/2}(\tilde{\boldsymbol{\theta}}_k - \tilde{\boldsymbol{\theta}}_{k'})$ we have

$$\begin{aligned} |\mathbf{e}_j^\top \tilde{\boldsymbol{\theta}}_k - \mathbf{e}_j^\top \tilde{\boldsymbol{\theta}}_{k'}|^2 &\leq \|\tilde{\boldsymbol{\theta}}_k - \tilde{\boldsymbol{\theta}}_{k'}\|^2 \\ &= \|\mathbf{B}_k^{-1/2} \mathbf{B}_k^{1/2}(\tilde{\boldsymbol{\theta}}_k - \tilde{\boldsymbol{\theta}}_{k'})\|^2 \\ &\leq \frac{\bar{\sigma}_{\max}^2(k)}{nh_k^d \Lambda_0} \|\mathbf{B}_k^{1/2}(\tilde{\boldsymbol{\theta}}_k - \tilde{\boldsymbol{\theta}}_{k'})\|^2. \end{aligned}$$

□

To obtain the “componentwise” oracle risk bounds we need to recheck the “propagation property”. Firstly, notice that the “propagation conditions” (2.17) on the choice the critical values $\mathfrak{z}_1, \dots, \mathfrak{z}_{K-1}$ imply the similar bounds for the components $\mathbf{e}_j^\top \hat{\boldsymbol{\theta}}_k(x)$. Recall that $\hat{\boldsymbol{\theta}}_k \stackrel{\text{def}}{=} \tilde{\boldsymbol{\theta}}_{\min\{k, \hat{k}\}}$. Then, by (2.17), Lemma 3.4, and the pivotality property from Lemma 4.1, we have the following simple observation:

Lemma 3.5. *Under the propagation conditions (PC) for any $\boldsymbol{\theta} \in \mathbb{R}^p$ and all $k = 2, \dots, K$ we have:*

$$\begin{aligned} \left(\frac{nh_k^d \Lambda_0}{\bar{\sigma}_{\max}^2(k)} \right)^r \mathbb{E}_{\boldsymbol{\theta}, \Sigma} |\mathbf{e}_j^\top \tilde{\boldsymbol{\theta}}_k(x) - \mathbf{e}_j^\top \hat{\boldsymbol{\theta}}_k(x)|^{2r} &\leq \mathbb{E}_{0, \Sigma} \|\mathbf{B}_k^{1/2}(\tilde{\boldsymbol{\theta}}_k - \hat{\boldsymbol{\theta}}_k)\|^{2r} \\ &\leq \alpha C(p, r). \end{aligned}$$

Here $\mathbb{E}_{0, \Sigma}$ stands for the expectation w.r.t. $\mathcal{N}(0, \Sigma)$ and $C(p, r)$ is given by (4.8).

As before we suppress the dependence on x . To get the propagation property we study for $k = 1, \dots, K$ the joint distributions of $\mathbf{e}_j^\top \tilde{\boldsymbol{\theta}}_1, \dots, \mathbf{e}_j^\top \tilde{\boldsymbol{\theta}}_k$, that is the distribution of $\mathbf{e}_j^\top \tilde{\boldsymbol{\Theta}}_k$, the j th row of the matrix $\tilde{\boldsymbol{\Theta}}_k$. Obviously,

$$\begin{aligned} \mathbb{E}_{\mathbf{f}, \Sigma_0} [\mathbf{e}_j^\top \tilde{\boldsymbol{\Theta}}_k] &= \mathbf{e}_j^\top \boldsymbol{\Theta}_k^* = (\mathbf{e}_j^\top \boldsymbol{\theta}_1^*, \dots, \mathbf{e}_j^\top \boldsymbol{\theta}_k^*), \\ \mathbb{E}_{\boldsymbol{\theta}, \Sigma} [\mathbf{e}_j^\top \tilde{\boldsymbol{\Theta}}_k] &= \mathbf{e}_j^\top \boldsymbol{\Theta}_k = (\mathbf{e}_j^\top \boldsymbol{\theta}, \dots, \mathbf{e}_j^\top \boldsymbol{\theta}). \end{aligned}$$

Recall that the matrices $\boldsymbol{\Sigma}_{k,0}$ and $\boldsymbol{\Sigma}_k$ have a block structure. Now, for instance, to study the estimator of the first coordinate of the vector $\boldsymbol{\theta} = \boldsymbol{\theta}(x)$, or of $f(x)$ in the case of the polynomial basis, we take the first elements of each block and so on. Denote the $k \times k$ covariance matrices of the j th elements of the vectors $\tilde{\boldsymbol{\theta}}_1, \dots, \tilde{\boldsymbol{\theta}}_k$ by

$$\begin{aligned} \boldsymbol{\Sigma}_{k,j} &\stackrel{\text{def}}{=} \left\{ \text{cov}_{\boldsymbol{\theta}, \Sigma} [\tilde{\theta}_l^{(j)}, \tilde{\theta}_m^{(j)}] \right\}_{1 \leq l \leq m \leq k} \\ &= \mathbf{D}_{k,j} (J_k \otimes \Sigma) \mathbf{D}_{k,j}^\top, \end{aligned} \tag{3.21}$$

$$\begin{aligned} \boldsymbol{\Sigma}_{k,0,j} &\stackrel{\text{def}}{=} \left\{ \text{cov}_{\mathbf{f}, \Sigma_0} [\tilde{\theta}_l^{(j)}, \tilde{\theta}_m^{(j)}] \right\}_{1 \leq l \leq m \leq k} \\ &= \mathbf{D}_{k,j} (J_k \otimes \Sigma_0) \mathbf{D}_{k,j}^\top, \end{aligned} \tag{3.22}$$

where J_k is a $k \times k$ matrix with all its elements equal to 1, and the $k \times nk$ block diagonal matrices $\mathbf{D}_{k,j}$ is defined by

$$\begin{aligned} \mathbf{D}_{k,j} &\stackrel{\text{def}}{=} \mathbf{e}_j^\top D_1 \oplus \dots \oplus \mathbf{e}_j^\top D_k = (I_k \otimes \mathbf{e}_j^\top) \mathbf{D}_k \\ D_l &\stackrel{\text{def}}{=} \mathbf{B}_l^{-1} \boldsymbol{\Psi} \mathbf{W}_l, \quad l = 1, \dots, k. \end{aligned} \tag{3.23}$$

Moreover, the following representation holds:

$$\begin{aligned}\Sigma_{k,j} &= (I_k \otimes e_j^\top) \mathbf{D}_k (J_k \otimes \Sigma) \mathbf{D}_k^\top (I_k \otimes e_j^\top)^\top \\ &= (I_k \otimes e_j)^\top \Sigma_k (I_k \otimes e_j),\end{aligned}\tag{3.24}$$

where Σ_k is defined by (3.5). Similarly,

$$\Sigma_{k,0,j} = (I_k \otimes e_j)^\top \Sigma_{k,0} (I_k \otimes e_j).\tag{3.25}$$

Thus, the important relation (3.8) is preserved for $\Sigma_{k,j}$ and $\Sigma_{k,0,j}$ obtained by picking up the (j,j) th elements of each block of Σ_k and $\Sigma_{k,0}$ respectively.

With usual notation $\gamma^{(j)}$ for the j th component of $\gamma \in \mathbb{R}^k$, denote by

$$\begin{aligned}b_j(k) &\stackrel{\text{def}}{=} (e_j^\top (\theta_1^* - \theta), \dots, e_j^\top (\theta_k^* - \theta))^\top \\ &= ((\theta_1^* - \theta)^{(j)}, \dots, (\theta_k^* - \theta)^{(j)})^\top \in \mathbb{R}^k\end{aligned}\tag{3.26}$$

$$\Delta_j(k) \stackrel{\text{def}}{=} b_j(k)^\top \Sigma_{k,j}^{-1} b_j(k).\tag{3.27}$$

Theorem 3.6. “Componentwise” propagation property

Under the conditions (A1) – (A7) and (PC) for any $k \leq K$ the following upper bound holds:

$$\begin{aligned}&\left(\frac{nh_k^d \Lambda_0}{\bar{\sigma}_{\max}^2(k)}\right)^{r/2} \mathbb{E}|e_j^\top \tilde{\theta}_k(x) - e_j^\top \hat{\theta}_k(x)|^r \\ &\leq (\alpha \mathbb{E}|\chi_p^2|^r)^{1/2} (1 + \delta)^{pk/4} (1 - \delta)^{-3pk/4} \exp\left\{\varphi(\delta) \frac{\Delta_j(k)}{2(1 - \delta)}\right\}\end{aligned}\tag{3.28}$$

with $\varphi(\delta)$ as in Theorem 3.2.

Corollary 3.7. Let the basis be polynomial. Then under the conditions of the preceding theorem $\mathbb{E}|\tilde{f}_k^{(j-1)}(x) - \hat{f}_k^{(j-1)}(x)|^r$ satisfy (3.28)

Proof. The proof essentially follows the line of the proof of Theorem 3.2. If the distributions of $\text{vec } \Theta_k$ were Gaussian, then any subvector is also Gaussian.

Denote by $\mathbb{P}_{\theta,\Sigma}^{k,j} = \mathcal{N}((e_j^\top \theta, \dots, e_j^\top \theta)^\top, \Sigma_{k,j})$ and by $\mathbb{P}_{f,\Sigma_0}^{k,j} = \mathcal{N}((e_j^\top \theta_1^*, \dots, e_j^\top \theta_k^*)^\top, \Sigma_{k,0,j})$, $k = 1, \dots, K$, the distributions of $e_j^\top \tilde{\theta}_k$ under the null and under the alternative.

By the Cauchy-Schwarz inequality and Lemma 3.5

$$\left(\frac{nh_k^d \Lambda_0}{\bar{\sigma}_{\max}^2(k)}\right)^{r/2} \mathbb{E}|e_j^\top \tilde{\theta}_k(x) - e_j^\top \hat{\theta}(x)|^r \leq (\alpha \mathbb{E}|\chi_p^2|^r)^{1/2} (\mathbb{E}_{\theta,\Sigma}[Z_{k,j}^2])^{1/2}$$

with the Radon-Nikodym derivative given by $Z_{k,j} = d\mathbb{P}_{f,\Sigma_0}^{k,j}/d\mathbb{P}_{\theta,\Sigma}^{k,j}$. By inequalities (3.24) and (3.25) the analog of Condition (A4) is preserved for $\Sigma_{k,0,j}$ and $\Sigma_{k,j}$, that is, there exist $\delta \in [0, 1)$ such that

$$(1 - \delta)\Sigma_{k,j} \preceq \Sigma_{k,0,j} \preceq (1 + \delta)\Sigma_{k,j}\tag{3.29}$$

for any $k \leq K$ and $j = 1, \dots, p$. Then the assertion of the theorem follows by the Taylor expansion at the point $(e_j^\top \theta, \dots, e_j^\top \theta)^\top$ and (3.29) similarly to the proof of Theorem 3.2. \square

At this point we introduce the “componentwise” small modeling bias conditions:

(SMBj) Let for some $j = 1, \dots, p$, some $k(j) \leq K$, and some $\theta^{(j)} = e_j^\top \boldsymbol{\theta}$ exist a constant $\Delta_j \geq 0$ such that

$$\Delta_j(k(j)) \leq \Delta_j, \quad (3.30)$$

where $\Delta_j(k)$ is defined by (3.27).

Definition 3.8. For each $j = 1, \dots, p$ the oracle index $k^*(j)$ is defined as the largest index in the scale for which the (SMBj) condition holds, that is

$$k^*(j) = \max\{k \leq K : \Delta_j(k) \leq \Delta_j\}. \quad (3.31)$$

Theorem 3.9. Assume (A1) – (A7) and (PC). Let the smallest bandwidth h_1 be such that the first estimate $e_j^\top \tilde{\boldsymbol{\theta}}_1(x)$ be always accepted in the adaptive procedure. Let $k^*(j)$ be the oracle index defined by (3.31), $j = 1, \dots, p$. Then the risk between the j th coordinates of the adaptive estimate and the oracle one is bounded with the following expression:

$$\begin{aligned} & \left(\frac{nh_{k^*(j)}^d \Lambda_0}{\bar{\sigma}_{\max}^2(k^*)} \right)^{r/2} \mathbb{E} |e_j^\top \tilde{\boldsymbol{\theta}}_{k^*(j)}(x) - e_j^\top \hat{\boldsymbol{\theta}}(x)|^r \\ & \leq \mathfrak{z}_{k^*(j)}^{r/2} + (\alpha \mathbb{E} |\chi_p^2|^r)^{1/2} (1 + \delta)^{pk_j^*/4} (1 - \delta)^{-3pk_j^*/4} \exp \left\{ \varphi(\delta) \frac{\Delta_j}{2(1 - \delta)} \right\}, \end{aligned} \quad (3.32)$$

where $\varphi(\delta)$ is as in Theorem 3.6.

Corollary 3.10. Let the basis be polynomial. Then under the conditions of the preceding theorem the risk between the adaptive estimate and the oracle one $\mathbb{E} |\tilde{f}_{k^*(j)}^{(j-1)}(x) - \hat{f}^{(j-1)}(x)|^r$ satisfy (3.32).

Proof. To simplify the notation we suppress the dependence on j in the index k . Similarly to the proof of Theorem 3.3 we consider the disjunct events $\{\hat{k} \leq k^*\}$ and $\{\hat{k} > k^*\}$. Therefore,

$$\begin{aligned} & \mathbb{E} |e_j^\top \tilde{\boldsymbol{\theta}}_{k^*}(x) - e_j^\top \hat{\boldsymbol{\theta}}(x)|^r \\ & = \mathbb{E} |e_j^\top \tilde{\boldsymbol{\theta}}_{k^*}(x) - e_j^\top \hat{\boldsymbol{\theta}}(x)|^r \mathbb{I}\{\hat{k} \leq k^*\} \\ & + \mathbb{E} |e_j^\top \tilde{\boldsymbol{\theta}}_{k^*}(x) - e_j^\top \hat{\boldsymbol{\theta}}(x)|^r \mathbb{I}\{\hat{k} > k^*\}. \end{aligned}$$

By Lemma 3.4 and the definition of the test statistic $T_{k^*, \hat{k}}$ the second summand can be easily bounded:

$$\begin{aligned} & \left(\frac{nh_{k^*}^d \Lambda_0}{\bar{\sigma}_{\max}^2(k^*)} \right)^{r/2} \mathbb{E} |e_j^\top \tilde{\boldsymbol{\theta}}_{k^*}(x) - e_j^\top \hat{\boldsymbol{\theta}}(x)|^r \mathbb{I}\{\hat{k} > k^*\} \\ & \leq \mathbb{E} \|\mathbf{B}_{k^*}^{1/2} (\tilde{\boldsymbol{\theta}}_{k^*}(x) - \hat{\boldsymbol{\theta}}(x))\|^r \mathbb{I}\{\hat{k} > k^*\} \\ & \leq \mathfrak{z}_{k^*}^{r/2}. \end{aligned}$$

To bound the first summand we use the “componentwise” analog of Theorem 3.2, particularly, Theorem 3.6, and this completes the proof. \square

3.5 SMB and the bias-variance trade-off

In [27] it was shown that the small modeling bias (*SMB1*) condition (3.30) can be obtained from the “bias-variance trade-off” relations. Notice that our set-up includes the set-up from [27] as a particular case. To prove that the similar relation holds in the present case we need the following definition. Given a point x and the method of localization w , for any $j = 1, \dots, p$ the “ideal adaptive bandwidths”, see [22], [23] is defined as follows:

$$k^*(j) = \max\{k \leq K : \bar{b}_{k,f^{(j-1)}}(x) \leq C_j(w)\sigma_k(x)\sqrt{d(n)}\}, \quad (3.33)$$

where $C_j(w)$ is a constant depending on the choice of the smoother w ,

$$\begin{aligned} \bar{b}_{k,f^{(j-1)}}(x) &= \sup_{1 \leq l \leq k} |\mathbf{e}_j^\top \boldsymbol{\theta}_l^*(x) - f^{(j-1)}(x)|, \\ \sigma_k^2(x) &= \text{Var}_{\mathbf{f}, \Sigma_0} [\mathbf{e}_j^\top \tilde{\boldsymbol{\theta}}_k(x)], \\ d(n) &= \log(h_K/h_1), \end{aligned}$$

and $f^{(0)}$ stands for the function f itself. To bound the “modeling bias” $\Delta_j(k)$ we need the following assumption:

(A8) *There exists a constant $s_j > 0$ such that for all $k \leq K$*

$$\boldsymbol{\Sigma}_{k,j}^{-1} \preceq s_j \boldsymbol{\Sigma}_{k,j,diag}^{-1} \quad (3.34)$$

where $\boldsymbol{\Sigma}_{k,j,diag} = \text{diag}(\text{Var}_{\boldsymbol{\theta}, \Sigma}[\mathbf{e}_j^\top \tilde{\boldsymbol{\theta}}_1(x)], \dots, \text{Var}_{\boldsymbol{\theta}, \Sigma}[\mathbf{e}_j^\top \tilde{\boldsymbol{\theta}}_k(x)])$ is a diagonal matrix composed of the diagonal elements of $\boldsymbol{\Sigma}_{k,j}$. Thus we have the following result:

Theorem 3.11. *Assume (A4), (A5) and (A8). Let the weights $\{w_{k,i}(x)\}$ satisfy (4.15). Then for any given point x , smoothing function w , and $j = 1, \dots, p$, the choice of $k(j) = k^*(j)$ defined by the relation (3.33) with $d(n) = 1$ implies the (*SMBj*) condition $\Delta_j(k(j)) \leq \Delta_j$ with the constant $\Delta_j = s_j C_j^2(w)(1 + \delta)(1 - u_0^{-1})^{-1}$.*

Proof. Consider the quantity $b_j(k)^\top \boldsymbol{\Sigma}_{k,j,diag}^{-1} b_j(k)$. Suppose that $\mathbf{e}_j^\top \boldsymbol{\theta}(x) = f^{(j-1)}(x)$. In view of relation (4.15) for the weights $\{w_{l,i}(x)\}$ the form of the matrix $\boldsymbol{\Sigma}_{k,j,diag}$ is particularly simple:

$$\boldsymbol{\Sigma}_{k,j,diag} = \text{diag}(\mathbf{e}_j^\top \mathbf{B}_1^{-1} \mathbf{e}_j, \dots, \mathbf{e}_j^\top \mathbf{B}_k^{-1} \mathbf{e}_j).$$

Then by (A5) and (3.2)

$$\begin{aligned} b_j(k)^\top \boldsymbol{\Sigma}_{k,j,diag}^{-1} b_j(k) &= \sum_{l=1}^k \frac{|\mathbf{e}_j^\top (\boldsymbol{\theta}_l^* - \boldsymbol{\theta})|^2}{\mathbf{e}_j^\top \mathbf{B}_l^{-1} \mathbf{e}_j} \\ &\leq (\bar{b}_{k,f^{(j-1)}}(x))^2 \sum_{l=1}^k \frac{1}{\mathbf{e}_j^\top \mathbf{B}_l^{-1} \mathbf{e}_j} \\ &\leq \frac{(\bar{b}_{k,f^{(j-1)}}(x))^2}{\mathbf{e}_j^\top \mathbf{B}_k^{-1} \mathbf{e}_j} \sum_{l=1}^k u_0^{-(k-l)} \\ &\leq \frac{(\bar{b}_{k,f^{(j-1)}}(x))^2 (1 + \delta)}{\sigma_k^2(x)(1 - u_0^{-1})}. \end{aligned}$$

By (3.33) with $d(n) = 1$, the choice of $k = k^*(j)$ implies $(\bar{b}_{k,f^{(j-1)}}(x))^2 \leq C_j^2(w)\sigma_k^2(x)$. Thus

$$b_j(k)^\top \Sigma_{k,j,diag}^{-1} b_j(k) \leq (1 + \delta)C_j^2(w)(1 - u_0^{-1})^{-1}$$

and

$$\Delta_j(k) = b_j(k)^\top \Sigma_{k,j}^{-1} b_j(k) \leq s_j C_j^2(w)(1 + \delta)(1 - u_0^{-1})^{-1}.$$

□

Remark 3.4. Using the standard technique it is easy to derive from the above result that for estimation of functions over Hölder classes the methodology proposed in [18] and [27] and generalized in the present paper delivers the minimax rate of convergence up to a logarithmic factor.

4 Appendix

4.1 Pivotality and local parametric risk bounds

Lemma 4.1. *Pivotality property*

Let (A3) hold. Let $\theta_1 = \dots = \theta_\varkappa = \theta$ for $\varkappa \leq K$. Then for any $k \leq \varkappa$ the risk associated with the adaptive estimate at every step of the procedure does not depend on the parameter θ :

$$\mathbb{E}_\theta |(\tilde{\theta}_k - \hat{\theta}_k)^\top \mathbf{B}_k (\tilde{\theta}_k - \hat{\theta}_k)|^r = \mathbb{E}_0 |(\tilde{\theta}_k - \hat{\theta}_k)^\top \mathbf{B}_k (\tilde{\theta}_k - \hat{\theta}_k)|^r,$$

where \mathbb{E}_0 denotes the expectation w.r.t. the centered measure $\mathcal{N}(0, \Sigma)$ or $\mathcal{N}(0, \Sigma_0)$.

Proof. After the first k steps $\hat{\theta}_k$ coincides with one of $\tilde{\theta}_m$, $m \leq k$, and this event takes place if for some $l \leq m$ the statistics $T_{l,m+1} > \mathfrak{z}_l$. In view of the decomposition (2.8) it holds

$$\begin{aligned} & \{T_{l,m+1} > \mathfrak{z}_l \text{ for some } l = 1, \dots, m | H_{m+1}\} \\ &= \left\{ (\tilde{\theta}_l - \tilde{\theta}_{m+1})^\top \mathbf{B}_l (\tilde{\theta}_l - \tilde{\theta}_{m+1}) > \mathfrak{z}_l \text{ for some } l = 1, \dots, m | H_{m+1} \right\} \\ &= \left\{ \left\| \mathbf{B}_l^{1/2} \left(\mathbf{B}_l^{-1} \Psi \mathbf{W}_l \Sigma_0^{1/2} \boldsymbol{\varepsilon} - \mathbf{B}_{m+1}^{-1} \Psi \mathbf{W}_{m+1} \Sigma_0^{1/2} \boldsymbol{\varepsilon} \right) \right\|^2 > \mathfrak{z}_l, \quad l \leq m \right\} \end{aligned}$$

with $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, I_n)$. The probability of this event does not depend on the shift θ , so without loss of generality θ can be taken equal to zero. The risk associated with $\hat{\theta}_k$ admits the following decomposition:

$$\mathbb{E}_\theta |(\tilde{\theta}_k - \hat{\theta}_k)^\top \mathbf{B}_k (\tilde{\theta}_k - \hat{\theta}_k)|^r = \sum_{m=1}^{k-1} \mathbb{E}_\theta |(\tilde{\theta}_k - \tilde{\theta}_m)^\top \mathbf{B}_k (\tilde{\theta}_k - \tilde{\theta}_m)|^r \mathbb{I}\{\hat{\theta}_k = \tilde{\theta}_m\}.$$

Under the conditions of the lemma for all $m < k$ the joint distribution of $(\tilde{\theta}_k - \tilde{\theta}_m)^\top \mathbf{B}_k (\tilde{\theta}_k - \tilde{\theta}_m)$ does not depend on θ by the same argumentation. □

To justify the statistical properties of the considered procedure we need the following simple observation. Let for any $\theta, \theta' \in \Theta$ the corresponding log-likelihood ratio $L(\mathbf{W}_k, \theta, \theta')$ be defined by (2.11). Then

$$2L(\mathbf{W}_k, \theta, \theta') = \|\mathbf{W}_k^{1/2}(\mathbf{Y} - \Psi^\top \theta')\|^2 - \|\mathbf{W}_k^{1/2}(\mathbf{Y} - \Psi^\top \theta)\|^2.$$

Theorem 4.2. Quadratic shape of the fitted log-likelihood

Let for every $k = 1, \dots, K$ the fitted log likelihood (FLL) be defined as follows:

$$L(\mathbf{W}_k, \tilde{\boldsymbol{\theta}}_k, \boldsymbol{\theta}') \stackrel{\text{def}}{=} \max_{\boldsymbol{\theta} \in \Theta} L(\mathbf{W}_k, \boldsymbol{\theta}, \boldsymbol{\theta}').$$

Then

$$2L(\mathbf{W}_k, \tilde{\boldsymbol{\theta}}_k, \boldsymbol{\theta}) = (\tilde{\boldsymbol{\theta}}_k - \boldsymbol{\theta})^\top \mathbf{B}_k (\tilde{\boldsymbol{\theta}}_k - \boldsymbol{\theta}). \quad (4.1)$$

Proof. Notice that $L(\mathbf{W}_k, \boldsymbol{\theta})$ defined by (2.3) is quadratic in $\boldsymbol{\theta}$. The assertion follows from the Taylor expansion of the second order at the point $\tilde{\boldsymbol{\theta}}_k$, because it is the point of maximum, and the second derivative is a constant matrix \mathbf{B}_k . \square

Let the matrix \mathbf{S} be defined as follows:

$$\mathbf{S} \stackrel{\text{def}}{=} \Sigma_0^{1/2} \mathbf{W}_k \boldsymbol{\Psi}^\top \mathbf{B}_k^{-1} \boldsymbol{\Psi} \mathbf{W}_k \Sigma_0^{1/2}. \quad (4.2)$$

Then for the distribution of $L(\mathbf{W}_k, \tilde{\boldsymbol{\theta}}_k, \boldsymbol{\theta}_k^*)$ one observes so-called ‘‘Wilks phenomenon’’, see [7], described by the following theorem:

Theorem 4.3. *Let the regression model be given by (1.1) and the parameter maximizing the expected local log-likelihood $\boldsymbol{\theta}_k^* = \boldsymbol{\theta}_k^*(x)$ be defined by (2.6). Then for any $k = 1, \dots, K$ the following equality in distribution takes place:*

$$2L(\mathbf{W}_k, \tilde{\boldsymbol{\theta}}_k, \boldsymbol{\theta}_k^*) \stackrel{d}{=} \lambda_1(\mathbf{S}) \bar{\varepsilon}_1^2 + \dots + \lambda_p(\mathbf{S}) \bar{\varepsilon}_p^2 \quad (4.3)$$

with $p = \text{rank}(\mathbf{B}_k) = \dim \Theta = p$. Here $\lambda_1(\mathbf{S}), \dots, \lambda_p(\mathbf{S})$ are the non-zero eigenvalues of the matrix \mathbf{S} , and $\bar{\varepsilon}_i$ are independent standard normal random variables.

Moreover, under (A4) the maximal eigenvalue $\lambda_{\max}(\mathbf{S}) \leq 1 + \delta$, and for any $\mathfrak{z} > 0$

$$\mathbb{P} \left\{ 2L(\mathbf{W}_k, \tilde{\boldsymbol{\theta}}_k, \boldsymbol{\theta}_k^*) \geq \mathfrak{z} \right\} \leq \mathbb{P} \left\{ \eta \geq \mathfrak{z}/(1 + \delta) \right\}, \quad (4.4)$$

where η is a random variable distributed according to the χ^2 law with p degrees of freedom.

Remark 4.1. Generally, if the matrix \mathbf{B}_k is degenerated in (4.3) the number of terms $p \leq \dim \Theta$.

Proof. By Theorem 4.2 and the decomposition (2.8) it holds that:

$$\begin{aligned} 2L(\mathbf{W}_k, \tilde{\boldsymbol{\theta}}_k, \boldsymbol{\theta}_k^*) &= (\tilde{\boldsymbol{\theta}}_k - \boldsymbol{\theta}_k^*)^\top \mathbf{B}_k (\tilde{\boldsymbol{\theta}}_k - \boldsymbol{\theta}_k^*) \\ &= (\mathbf{B}_k^{-1} \boldsymbol{\Psi} \mathbf{W}_k \Sigma_0^{1/2} \boldsymbol{\varepsilon})^\top \mathbf{B}_k (\mathbf{B}_k^{-1} \boldsymbol{\Psi} \mathbf{W}_k \Sigma_0^{1/2} \boldsymbol{\varepsilon}) \\ &= \boldsymbol{\varepsilon}^\top \mathbf{S} \boldsymbol{\varepsilon}, \end{aligned}$$

where the symmetric matrix \mathbf{S} is defined by (4.2). Then by the Schur theorem there exist an orthogonal matrix \mathbf{M} and the diagonal matrix $\boldsymbol{\Lambda}$ composed of the eigenvalues of \mathbf{S} such that $\mathbf{S} = \mathbf{M}^\top \boldsymbol{\Lambda} \mathbf{M}$. For $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, I_n)$ and an orthogonal matrix \mathbf{M} it holds that $\bar{\boldsymbol{\varepsilon}} \stackrel{\text{def}}{=} \mathbf{M} \boldsymbol{\varepsilon} \sim \mathcal{N}(0, I_n)$. Indeed, $\mathbb{E} \mathbf{M} \boldsymbol{\varepsilon} = \mathbb{E} \boldsymbol{\varepsilon} = 0$ and

$$\text{Var } \mathbf{M} \boldsymbol{\varepsilon} = \mathbb{E} \mathbf{M} \boldsymbol{\varepsilon} (\mathbf{M} \boldsymbol{\varepsilon})^\top = \mathbf{M} \mathbb{E} (\boldsymbol{\varepsilon} \boldsymbol{\varepsilon}^\top) \mathbf{M} \boldsymbol{\varepsilon} = \mathbf{M} \mathbf{M}^\top = I_n.$$

Therefore,

$$2L(\mathbf{W}_k, \tilde{\boldsymbol{\theta}}_k, \boldsymbol{\theta}_k^*) \stackrel{d}{=} \bar{\boldsymbol{\varepsilon}}^\top \boldsymbol{\Lambda} \bar{\boldsymbol{\varepsilon}}, \quad \bar{\boldsymbol{\varepsilon}} \sim \mathcal{N}(0, I_n).$$

On the other hand, the matrix $\mathbf{S} = \Sigma_0^{1/2} \mathbf{W}_k \boldsymbol{\Psi}^\top \mathbf{B}_k^{-1} \boldsymbol{\Psi} \mathbf{W}_k \Sigma_0^{1/2}$ can be rewritten as:

$$\mathbf{S} = \Sigma_0^{1/2} \mathbf{W}_k^{1/2} \boldsymbol{\Pi}_k \mathbf{W}_k^{1/2} \Sigma_0^{1/2},$$

with $\boldsymbol{\Pi}_k = \mathbf{W}_k^{1/2} \boldsymbol{\Psi}^\top \mathbf{B}_k^{-1} \boldsymbol{\Psi} \mathbf{W}_k^{1/2}$. Notice that $\boldsymbol{\Pi}_k$ is an orthogonal projector on the linear subspace of dimension $p = \text{rank}(\mathbf{B}_k)$ spanned by the rows of matrix $\boldsymbol{\Psi}$. Indeed, $\boldsymbol{\Pi}_k$ is symmetric and idempotent, i.e., $\boldsymbol{\Pi}_k^2 = \boldsymbol{\Pi}_k$.

Moreover, $\text{rank}(\boldsymbol{\Pi}_k) = \text{tr}(\boldsymbol{\Pi}_k) = \text{tr}(\mathbf{W}_k^{1/2} \boldsymbol{\Psi}^\top \mathbf{B}_k^{-1} \boldsymbol{\Psi} \mathbf{W}_k^{1/2}) = \text{tr}(\mathbf{B}_k^{-1} \boldsymbol{\Psi} \mathbf{W}_k \boldsymbol{\Psi}^\top) = \text{tr}(\mathbf{B}_k^{-1} \mathbf{B}_k) = \text{tr}(I_p) = p$. Therefore $\boldsymbol{\Pi}_k$ has only p unit eigenvalues and $n-p$ zero ones. Notice also that the $n \times n$ matrix \mathbf{S} has $\text{rank}(\mathbf{S}) = \text{rank}(\boldsymbol{\Pi}_k \mathbf{W}_k^{1/2} \Sigma_0^{1/2}) = \text{rank}(\boldsymbol{\Pi}_k) = p$ as well. Thus $2L(\mathbf{W}_k, \tilde{\boldsymbol{\theta}}_k, \boldsymbol{\theta}_k^*) \stackrel{d}{=} \lambda_1(\mathbf{S}) \bar{\boldsymbol{\varepsilon}}_1^2 + \dots + \lambda_p(\mathbf{S}) \bar{\boldsymbol{\varepsilon}}_p^2$, where $\lambda_1(\mathbf{S}), \dots, \lambda_p(\mathbf{S})$ are the non-zero eigenvalues of the matrix \mathbf{S} .

Recall the definition of the matrix norm induced by the L_2 vector norm:

$$\|A\| \stackrel{\text{def}}{=} \sqrt{\lambda_{\max}(A^\top A)}. \quad (4.5)$$

Thus, taking into account Assumption (A4), the induced L_2 -norm of matrix \mathbf{S} can be estimated as follows:

$$\begin{aligned} \|\mathbf{S}\| &= \|\Sigma_0^{1/2} \mathbf{W}_k^{1/2} \boldsymbol{\Pi}_k \mathbf{W}_k^{1/2} \Sigma_0^{1/2}\| \\ &\leq \|\Sigma_0^{1/2} \mathbf{W}_k^{1/2}\| \|\boldsymbol{\Pi}_k\| \|\mathbf{W}_k^{1/2} \Sigma_0^{1/2}\| \\ &= \lambda_{\max}(\mathbf{W}_k \Sigma_0) \lambda_{\max}(\boldsymbol{\Pi}_k) \\ &= \max_i \left\{ w_{k,i} \frac{\sigma_{0,i}^2}{\sigma_i^2} \right\} \\ &\leq (1 + \delta) \max_i \{w_{k,i}\} \leq 1 + \delta. \end{aligned}$$

Therefore, the largest eigenvalue of matrix \mathbf{S} is bounded: $\lambda_{\max}(\mathbf{S}) \leq 1 + \delta$.

The last assertion of the theorem follows from the simple observation that

$$I\!\!P \{ \lambda_1(\mathbf{S}) \bar{\boldsymbol{\varepsilon}}_1^2 + \dots + \lambda_p(\mathbf{S}) \bar{\boldsymbol{\varepsilon}}_p^2 \geq \mathfrak{z} \} \leq I\!\!P \{ \lambda_{\max}(\mathbf{S}) (\bar{\boldsymbol{\varepsilon}}_1^2 + \dots + \bar{\boldsymbol{\varepsilon}}_p^2) \geq \mathfrak{z} \}.$$

□

Corollary 4.4. Quasi-parametric risk bounds

Let the model be given by (1.1) and $\boldsymbol{\theta}_k^* = \boldsymbol{\theta}_k^*(x)$ be defined by (2.6). Assume (A4). Then for any $\mu < 1/(1 + \delta)$

$$\mathbb{E} \exp\{\mu L(\mathbf{W}_k, \tilde{\boldsymbol{\theta}}_k, \boldsymbol{\theta}_k^*)\} \leq [1 - \mu(1 + \delta)]^{-p/2} \quad (4.6)$$

$$\mathbb{E}|2L(\mathbf{W}_k, \tilde{\boldsymbol{\theta}}_k, \boldsymbol{\theta}_k^*)|^r \leq (1 + \delta)^r C(p, r), \quad (4.7)$$

where

$$C(p, r) = \mathbb{E}|\chi_p^2|^r = 2^r \Gamma(r + p/2)/\Gamma(p/2). \quad (4.8)$$

Proof. By (4.3) and the independence of $\bar{\varepsilon}_i$

$$\begin{aligned}
\mathbb{E} \exp\{\mu L(\mathbf{W}_k, \tilde{\boldsymbol{\theta}}_k, \boldsymbol{\theta}_k^*)\} &= \mathbb{E} \exp\left\{\frac{\mu}{2} \sum_{i=1}^p \lambda_i(\mathbf{S}) \bar{\varepsilon}_i^2\right\} \\
&= \prod_{i=1}^p \mathbb{E} \exp\left\{\frac{\mu}{2} \lambda_i(\mathbf{S}) \bar{\varepsilon}_i^2\right\} \\
&= \prod_{i=1}^p [1 - \mu \lambda_i(\mathbf{S})]^{-1/2} \\
&\leq [1 - \mu \lambda_{max}(\mathbf{S})]^{-p/2} \\
&\leq [1 - \mu(1 + \delta)]^{-p/2}.
\end{aligned}$$

Let $\eta \sim \chi_p^2$. Integrating by parts yields the second inequality:

$$\begin{aligned}
\mathbb{E}|2L(\mathbf{W}_k, \tilde{\boldsymbol{\theta}}_k, \boldsymbol{\theta}_k^*)|^r &= \int_0^\infty \mathbb{I}\{2L(\mathbf{W}_k, \tilde{\boldsymbol{\theta}}_k, \boldsymbol{\theta}_k^*) \geq \mathfrak{z}\} r \mathfrak{z}^{r-1} d\mathfrak{z} \\
&\leq r \int_0^\infty \mathbb{I}\{\eta \geq \mathfrak{z}/(1 + \delta)\} \mathfrak{z}^{r-1} d\mathfrak{z} \\
&= (1 + \delta)^r \mathbb{E}|\eta|^r.
\end{aligned}$$

□

4.2 Proof of the bounds for the critical values

Denote for any $l < k$ the variance of the difference $\tilde{\boldsymbol{\theta}}_k - \tilde{\boldsymbol{\theta}}_l$ by V_{lk} :

$$V_{lk} \stackrel{\text{def}}{=} \text{Var}(\tilde{\boldsymbol{\theta}}_k - \tilde{\boldsymbol{\theta}}_l) \succ 0. \quad (4.9)$$

Then there exists a unique matrix $V_{lk}^{1/2} \succ 0$ such that $(V_{lk}^{1/2})^2 = V_{lk}$.

Lemma 4.5. *Assume (A4), (A3) and (A5). If for some $k \leq K$ the hypothesis H_k is true, that is, if $\boldsymbol{\theta}_1^* = \dots = \boldsymbol{\theta}_k^* = \boldsymbol{\theta}$, then for any $l < k$ it holds that:*

$$\begin{aligned}
\mathbb{I}\{2L(\mathbf{W}_l, \tilde{\boldsymbol{\theta}}_l, \tilde{\boldsymbol{\theta}}_k) \geq \mathfrak{z}\} &\leq \mathbb{I}\{\eta \geq \mathfrak{z}/\lambda_{max}(V_{lk}^{1/2} \mathbf{B}_l V_{lk}^{1/2})\} \\
&\leq \mathbb{I}\{\eta \geq \mathfrak{z}/t_0\} \\
\mathbb{I}\{2L(\mathbf{W}_k, \tilde{\boldsymbol{\theta}}_k, \tilde{\boldsymbol{\theta}}_l) \geq \mathfrak{z}\} &\leq \mathbb{I}\{\eta \geq \mathfrak{z}/\lambda_{max}(V_{lk}^{1/2} \mathbf{B}_k V_{lk}^{1/2})\} \\
&\leq \mathbb{I}\{\eta \geq \mathfrak{z}/t_1\}
\end{aligned}$$

where $t_0 = 2(1 + \delta)(1 + u_0^{-(k-l)})$, $t_1 = 2(1 + \delta)(1 + u^{(k-l)})$, and η is the χ_p^2 -distributed random variable.

Proof. The H_k and (2.8) imply

$$\tilde{\boldsymbol{\theta}}_l - \tilde{\boldsymbol{\theta}}_k = \mathbf{B}_l^{-1} \boldsymbol{\Psi} \mathbf{W}_l \Sigma_0^{1/2} \boldsymbol{\varepsilon} - \mathbf{B}_k^{-1} \boldsymbol{\Psi} \mathbf{W}_k \Sigma_0^{1/2} \boldsymbol{\varepsilon} \stackrel{d}{=} V_{lk}^{1/2} \boldsymbol{\xi},$$

where ξ is a standard normal vector in \mathbb{R}^p . Thus by Theorem 4.2 for any $l < k$

$$2L(\mathbf{W}_l, \tilde{\boldsymbol{\theta}}_l, \tilde{\boldsymbol{\theta}}_k) = \|\mathbf{B}_l^{1/2}(\tilde{\boldsymbol{\theta}}_l - \tilde{\boldsymbol{\theta}}_k)\|^2 \stackrel{d}{=} \xi^\top V_{lk}^{1/2} \mathbf{B}_l V_{lk}^{1/2} \xi.$$

By the Schur theorem there exists an orthogonal matrix M such that

$$\xi^\top V_{lk}^{1/2} \mathbf{B}_l V_{lk}^{1/2} \xi \stackrel{d}{=} \bar{\varepsilon}^\top M^\top \Lambda_{lk} M \bar{\varepsilon},$$

where $\bar{\varepsilon}$ is a standard normal vector, $\Lambda = \text{diag}(\lambda_1(V_{lk}^{1/2} \mathbf{B}_l V_{lk}^{1/2})), \dots, \lambda_p(V_{lk}^{1/2} \mathbf{B}_l V_{lk}^{1/2}))$, and $p = \text{rank}(\mathbf{B}_l)$. Therefore,

$$2L(\mathbf{W}_l, \tilde{\boldsymbol{\theta}}_l, \tilde{\boldsymbol{\theta}}_k) \stackrel{d}{=} \lambda_1(V_{lk}^{1/2} \mathbf{B}_l V_{lk}^{1/2}) \bar{\varepsilon}_1^2 + \dots + \lambda_p(V_{lk}^{1/2} \mathbf{B}_l V_{lk}^{1/2}) \bar{\varepsilon}_p^2,$$

where $\lambda_j(V_{lk}^{1/2} \mathbf{B}_l V_{lk}^{1/2})$, $j = 1, \dots, p$, are the nonzero eigenvalues of $V_{lk}^{1/2} \mathbf{B}_l V_{lk}^{1/2}$.

By the similar argumentation:

$$2L(\mathbf{W}_k, \tilde{\boldsymbol{\theta}}_k, \tilde{\boldsymbol{\theta}}_l) \stackrel{d}{=} \lambda_1(V_{lk}^{1/2} \mathbf{B}_k V_{lk}^{1/2}) \bar{\varepsilon}_1^2 + \dots + \lambda_p(V_{lk}^{1/2} \mathbf{B}_k V_{lk}^{1/2}) \bar{\varepsilon}_p^2.$$

Denote by η the χ_p^2 -distributed random variable, then

$$\begin{aligned} I\!\!P \left\{ 2L(\mathbf{W}_l, \tilde{\boldsymbol{\theta}}_l, \tilde{\boldsymbol{\theta}}_k) \geq \mathfrak{z} \right\} &\leq I\!\!P \left\{ \eta \geq \mathfrak{z}/\lambda_{\max}(V_{lk}^{1/2} \mathbf{B}_l V_{lk}^{1/2}) \right\} \\ I\!\!P \left\{ 2L(\mathbf{W}_k, \tilde{\boldsymbol{\theta}}_k, \tilde{\boldsymbol{\theta}}_l) \geq \mathfrak{z} \right\} &\leq I\!\!P \left\{ \eta \geq \mathfrak{z}/\lambda_{\max}(V_{lk}^{1/2} \mathbf{B}_k V_{lk}^{1/2}) \right\} \end{aligned}$$

For any square matrices A and B we have $(A - B)(A^\top - B^\top) \preceq 2(AA^\top + BB^\top)$. Application of this bound to the variance of the difference of estimates yields

$$\begin{aligned} V_{lk} &= \left(\mathbf{B}_l^{-1} \Psi \mathbf{W}_l \Sigma_0^{1/2} - \mathbf{B}_k^{-1} \Psi \mathbf{W}_k \Sigma_0^{1/2} \right) \left(\mathbf{B}_l^{-1} \Psi \mathbf{W}_l \Sigma_0^{1/2} - \mathbf{B}_k^{-1} \Psi \mathbf{W}_k \Sigma_0^{1/2} \right)^\top \\ &\leq 2(\mathbf{B}_l^{-1} \Psi \mathbf{W}_l \Sigma_0 \mathbf{W}_l \Psi^\top \mathbf{B}_l^{-1} + \mathbf{B}_k^{-1} \Psi \mathbf{W}_k \Sigma_0 \mathbf{W}_k \Psi^\top \mathbf{B}_k^{-1}) \\ &= 2V_l + 2V_k, \end{aligned}$$

where $V_l = \text{Var } \tilde{\boldsymbol{\theta}}_l$, $l \leq k$. By (3.2) and by Assumption (B) we have:

$$\begin{aligned} V_l &\preceq (1 + \delta) \mathbf{B}_l^{-1}, \\ V_k &\preceq (1 + \delta) \mathbf{B}_k^{-1} \preceq (1 + \delta) u_0^{-(k-l)} \mathbf{B}_l^{-1}, \\ V_{lk} &\preceq 2(1 + \delta)(1 + u_0^{-(k-l)}) \mathbf{B}_l^{-1}. \end{aligned}$$

Therefore,

$$\mathbf{B}_l \preceq 2(1 + \delta)(1 + u_0^{-(k-l)}) V_{lk}^{-1}. \quad (4.10)$$

Thus by (4.10) the upper bound for the induced matrix norm reads as follows:

$$\begin{aligned} \lambda_{\max}(V_{lk}^{1/2} \mathbf{B}_l V_{lk}^{1/2}) &= \|\mathbf{B}_l^{1/2} V_{lk}^{1/2}\|^2 \\ &= \sup_{\|\gamma\|=1} \gamma^\top V_{lk}^{1/2} \mathbf{B}_l V_{lk}^{1/2} \gamma \\ &\leq 2(1 + \delta)(1 + u_0^{-(k-l)}) \sup_{\|\gamma\|=1} \gamma^\top V_{lk}^{1/2} V_{lk}^{-1} V_{lk}^{1/2} \gamma \\ &\leq 2(1 + \delta)(1 + u_0^{-(k-l)}). \end{aligned} \quad (4.11)$$

Similarly,

$$\begin{aligned} V_{lk} &\leq 2(1+\delta)(1+u^{(k-l)})\mathbf{B}_k^{-1}, \\ \lambda_{max}(V_{lk}^{1/2}\mathbf{B}_k V_{lk}^{1/2}) &\leq 2(1+\delta)(1+u^{(k-l)}). \end{aligned} \quad (4.12)$$

These bounds imply

$$\begin{aligned} \mathbb{P}\left\{2L(\mathbf{W}_l, \tilde{\boldsymbol{\theta}}_l, \tilde{\boldsymbol{\theta}}_k) \geq \mathfrak{z}\right\} &\leq \mathbb{P}\left\{\eta \geq \mathfrak{z}/\lambda_{max}(V_{lk}^{1/2}\mathbf{B}_l V_{lk}^{1/2})\right\} \\ &\leq \mathbb{P}\left\{\eta \geq \mathfrak{z}\left[2(1+\delta)(1+u_0^{-(k-l)})\right]^{-1}\right\} \\ \mathbb{P}\left\{2L(\mathbf{W}_k, \tilde{\boldsymbol{\theta}}_k, \tilde{\boldsymbol{\theta}}_l) \geq \mathfrak{z}\right\} &\leq \mathbb{P}\left\{\eta \geq \mathfrak{z}/\lambda_{max}(V_{lk}^{1/2}\mathbf{B}_k V_{lk}^{1/2})\right\} \\ &\leq \mathbb{P}\left\{\eta \geq \mathfrak{z}\left[2(1+\delta)(1+u^{(k-l)})\right]^{-1}\right\} \end{aligned}$$

□

Lemma 4.6. Under the conditions of the preceding lemma for any $\mu_0 < t_0^{-1}$, or $\mu_1 < t_1^{-1}$ respectively, the exponential moments are bounded:

$$\begin{aligned} \mathbb{E} \exp\{\mu_0 L(\mathbf{W}_l, \tilde{\boldsymbol{\theta}}_l, \tilde{\boldsymbol{\theta}}_k)\} &\leq [1 - \mu_0 t_0]^{-p/2} \\ \mathbb{E} \exp\{\mu_1 L(\mathbf{W}_k, \tilde{\boldsymbol{\theta}}_k, \tilde{\boldsymbol{\theta}}_l)\} &\leq [1 - \mu_1 t_1]^{-p/2}, \end{aligned}$$

where $t_0 = 2(1+\delta)(1+u_0^{-(k-l)})$ and $t_1 = 2(1+\delta)(1+u^{(k-l)})$.

Proof. The statement of the lemma is justified similarly to the proof of Corollary 4.4. The bounds (4.11) and (4.12) imply the bounds for the corresponding moment generating functions:

$$\begin{aligned} \mathbb{E} \exp\{\mu L(\mathbf{W}_l, \tilde{\boldsymbol{\theta}}_l, \tilde{\boldsymbol{\theta}}_k)\} &= \prod_{j=1}^p \mathbb{E} \exp\left\{\frac{\mu}{2} \lambda_j (V_{lk}^{1/2} \mathbf{B}_l V_{lk}^{1/2}) \bar{\varepsilon}_j^2\right\} \\ &= \prod_{j=1}^p [1 - \mu \lambda_j (V_{lk}^{1/2} \mathbf{B}_l V_{lk}^{1/2})]^{-1/2} \\ &\leq [1 - \mu \lambda_{max}(V_{lk}^{1/2} \mathbf{B}_l V_{lk}^{1/2})]^{-p/2} \\ &\leq [1 - 2\mu(1+\delta)(1+u_0^{-(k-l)})]^{-p/2}, \\ \mathbb{E} \exp\{\mu L(\mathbf{W}_k, \tilde{\boldsymbol{\theta}}_k, \tilde{\boldsymbol{\theta}}_l)\} &\leq [1 - \mu \lambda_{max}(V_{lk}^{1/2} \mathbf{B}_k V_{lk}^{1/2})]^{-p/2} \\ &\leq [1 - 2\mu(1+\delta)(1+u^{(k-l)})]^{-p/2}. \end{aligned}$$

□

Lemma 4.7. Under the conditions of the preceding lemma it holds that:

$$\begin{aligned} \mathbb{E}|2L(\mathbf{W}_l, \tilde{\boldsymbol{\theta}}_l, \tilde{\boldsymbol{\theta}}_k)|^r &\leq 2^r C(p, r)(1+\delta)^r (1+u_0^{-(k-l)})^r, \\ \mathbb{E}|2L(\mathbf{W}_k, \tilde{\boldsymbol{\theta}}_k, \tilde{\boldsymbol{\theta}}_l)|^r &\leq 2^r C(p, r)(1+\delta)^r (1+u^{(k-l)})^r, \end{aligned}$$

where $C(p, r)$ is as in (4.8).

Proof. Integrating by parts and Lemma 4.5 yield for the second assertion

$$\begin{aligned}\mathbb{E}|2L(\mathbf{W}_k, \tilde{\boldsymbol{\theta}}_k, \tilde{\boldsymbol{\theta}}_l)|^r &= r \int_0^\infty \mathbb{P} \left\{ 2L(\mathbf{W}_k, \tilde{\boldsymbol{\theta}}_k, \tilde{\boldsymbol{\theta}}_l) \geq \mathfrak{z} \right\} \mathfrak{z}^{r-1} d\mathfrak{z} \\ &\leq r \int_0^\infty \mathbb{P} \left\{ \eta \geq \mathfrak{z} \left[2(1+\delta)(1+u^{(k-l)}) \right]^{-1} \right\} \mathfrak{z}^{r-1} d\mathfrak{z} \\ &= 2^r (1+\delta)^r (1+u^{(k-l)})^r \mathbb{E}|\eta|^r,\end{aligned}$$

where $\eta \sim \chi_p^2$. The first assertion is proved similarly. \square

Proof. of Theorem 3.1 The theoretical choice of the critical values The risk corresponding to the adaptive estimate can be represented as a sum of risks of the false alarms at each step of the procedure:

$$\mathbb{E}_{0,\Sigma} |(\tilde{\boldsymbol{\theta}}_k - \hat{\boldsymbol{\theta}}_k)^\top \mathbf{B}_k (\tilde{\boldsymbol{\theta}}_k - \hat{\boldsymbol{\theta}}_k)|^r = \sum_{m=1}^{k-1} \mathbb{E}_{0,\Sigma} |(\tilde{\boldsymbol{\theta}}_k - \tilde{\boldsymbol{\theta}}_m)^\top \mathbf{B}_k (\tilde{\boldsymbol{\theta}}_k - \tilde{\boldsymbol{\theta}}_m)|^r \mathbb{I}\{\hat{\boldsymbol{\theta}}_k = \tilde{\boldsymbol{\theta}}_m\}.$$

By the definition of the last accepted estimate $\hat{\boldsymbol{\theta}}_k$, for any $m = 1, \dots, k-1$, the event $\{\hat{\boldsymbol{\theta}}_k = \tilde{\boldsymbol{\theta}}_m\}$ happens if for some $l = 1, \dots, m$ the statistic $T_{l,m+1} > \mathfrak{z}_l$. Thus

$$\{\hat{\boldsymbol{\theta}}_k = \tilde{\boldsymbol{\theta}}_m\} \subseteq \bigcup_{l=1}^m \{T_{l,m+1} > \mathfrak{z}_l\}.$$

It holds also that for any positive μ

$$\begin{aligned}\mathbb{I}\{T_{l,m+1} > \mathfrak{z}_l\} &= \mathbb{I}\{2L(\mathbf{W}_l, \tilde{\boldsymbol{\theta}}_l, \tilde{\boldsymbol{\theta}}_{m+1}) - \mathfrak{z}_l > 0\} \\ &\leq \exp\left\{\frac{\mu}{2} L(\mathbf{W}_l, \tilde{\boldsymbol{\theta}}_l, \tilde{\boldsymbol{\theta}}_{m+1}) - \frac{\mu}{4} \mathfrak{z}_l\right\}.\end{aligned}$$

This simple fact and the Cauchy-Schwarz inequality imply for $m = 1, \dots, k-1$ the following bound:

$$\begin{aligned}&\mathbb{E}_{0,\Sigma} |(\tilde{\boldsymbol{\theta}}_k - \tilde{\boldsymbol{\theta}}_m)^\top \mathbf{B}_k (\tilde{\boldsymbol{\theta}}_k - \tilde{\boldsymbol{\theta}}_m)|^r \mathbb{I}\{\hat{\boldsymbol{\theta}}_k = \tilde{\boldsymbol{\theta}}_m\} \\ &= \mathbb{E}_{0,\Sigma} |2L(\mathbf{W}_k, \tilde{\boldsymbol{\theta}}_k, \tilde{\boldsymbol{\theta}}_m)|^r \mathbb{I}\{\hat{\boldsymbol{\theta}}_k = \tilde{\boldsymbol{\theta}}_m\} \\ &\leq \sum_{l=1}^m e^{-\frac{\mu}{4} \mathfrak{z}_l} \mathbb{E}_{0,\Sigma} \left[|2L(\mathbf{W}_k, \tilde{\boldsymbol{\theta}}_k, \tilde{\boldsymbol{\theta}}_m)|^r \exp\left\{\frac{\mu}{2} L(\mathbf{W}_l, \tilde{\boldsymbol{\theta}}_l, \tilde{\boldsymbol{\theta}}_{m+1})\right\} \right] \\ &\leq \sum_{l=1}^m e^{-\frac{\mu}{4} \mathfrak{z}_l} \left\{ \mathbb{E}_{0,\Sigma} \left[|2L(\mathbf{W}_k, \tilde{\boldsymbol{\theta}}_k, \tilde{\boldsymbol{\theta}}_m)|^{2r} \right] \right\}^{1/2} \left\{ \mathbb{E}_{0,\Sigma} \left[\exp\left\{\mu L(\mathbf{W}_l, \tilde{\boldsymbol{\theta}}_l, \tilde{\boldsymbol{\theta}}_{m+1})\right\} \right] \right\}^{1/2}.\end{aligned}$$

By Lemma 4.6 with $\delta = 0$

$$\mathbb{E}_{0,\Sigma} \left[\exp\left\{\mu L(\mathbf{W}_l, \tilde{\boldsymbol{\theta}}_l, \tilde{\boldsymbol{\theta}}_{m+1})\right\} \right] < (1-4\mu)^{-p/2}.$$

This together with the bound from Lemma 4.7 gives

$$\begin{aligned}
& \mathbb{E}_{0,\Sigma} |(\tilde{\boldsymbol{\theta}}_k - \hat{\boldsymbol{\theta}}_k)^\top \mathbf{B}_k (\tilde{\boldsymbol{\theta}}_k - \hat{\boldsymbol{\theta}}_k)|^r \\
& \leq 2^r \sqrt{C(p, 2r)} (1 - 4\mu)^{-p/4} \sum_{m=1}^{k-1} \sum_{l=1}^m e^{-\frac{\mu}{4}\mathfrak{z}_l} (1 + u^{(k-m)})^r \\
& = 2^r \sqrt{C(p, 2r)} (1 - 4\mu)^{-p/4} \sum_{l=1}^{k-1} e^{-\frac{\mu}{4}\mathfrak{z}_l} \sum_{m=l}^{k-1} (1 + u^{(k-m)})^r \\
& \leq 2^{2r} \sqrt{C(p, 2r)} (1 - 4\mu)^{-p/4} (1 - u^{-r})^{-1} \sum_{l=1}^{k-1} e^{-\frac{\mu}{4}\mathfrak{z}_l} u^{r(k-l)},
\end{aligned}$$

because $-(k-l) < -(m-l)$ and

$$\begin{aligned}
\sum_{m=l}^{k-1} (1 + u^{(k-m)})^r &= u^{r(k-l)} \sum_{m=l}^{k-1} (u^{-(k-l)} + u^{-(m-l)})^r \\
&< 2^r u^{r(k-l)} \sum_{m=l}^{k-1} u^{-r(m-l)} \\
&< 2^r u^{r(k-l)} (1 - u^{-r})^{-1}.
\end{aligned}$$

Since $u^{r(k-l)} \leq u^{r(K-l)}$ for any $l < k \leq K$ the choice

$$\mathfrak{z}_l = \frac{4}{\mu} \left\{ r(K-l) \log u + \log(K/\alpha) - \frac{p}{4} \log(1-4\mu) - \log(1-u^{-r}) + \bar{C}(p, r) \right\}.$$

with

$$\bar{C}(p, r) = \log \left\{ \frac{2^{2r} [\Gamma(2r+p/2)\Gamma(p/2)]^{1/2}}{\Gamma(r+p/2)} \right\}$$

provides the required bound

$$\mathbb{E}_{0,\Sigma} |(\tilde{\boldsymbol{\theta}}_l - \hat{\boldsymbol{\theta}}_l)^\top \mathbf{B}_l (\tilde{\boldsymbol{\theta}}_l - \hat{\boldsymbol{\theta}}_l)|^r \leq \alpha C(p, r) \quad \text{for all } l = 2, \dots, K.$$

□

4.3 Matrix results

Lemma 4.8. *The matrices $J_k \otimes \Sigma$ and $J_k \otimes \Sigma_0$ are positive semidefinite for any $k = 2, \dots, K$.*

Moreover, under Assumption (A4) with the same δ , the similar to (A4) relation holds for the covariance matrices $\boldsymbol{\Sigma}_k$ and $\boldsymbol{\Sigma}_{k,0}$ of linear estimates:

$$(1 - \delta) \boldsymbol{\Sigma}_k \preceq \boldsymbol{\Sigma}_{k,0} \preceq (1 + \delta) \boldsymbol{\Sigma}_k, \quad k \leq K.$$

Proof. Symmetry of J_k and Σ , (respectively, Σ_0) implies symmetry of $J_k \otimes \Sigma$, (respectively, $J_k \otimes \Sigma_0$). Notice that any vector $\gamma_{nk} \in \mathbb{R}^{nk}$ can be represented as a partitioned vector $\gamma_{nk}^\top = ((\gamma_{nk}^{(1)})^\top, (\gamma_{nk}^{(2)})^\top, \dots, (\gamma_{nk}^{(k)})^\top)$, with $\gamma_{nk}^{(l)} \in \mathbb{R}^n$, $l = 1, \dots, k$. Then

$$\gamma_{nk}^\top (J_k \otimes \Sigma) \gamma_{nk} = \left(\sum_{l=1}^k \gamma_{nk}^{(l)} \right)^\top \Sigma \left(\sum_{l=1}^k \gamma_{nk}^{(l)} \right) = \tilde{\gamma}_n^\top \Sigma \tilde{\gamma}_n, \quad (4.13)$$

where $\tilde{\gamma}_n \stackrel{\text{def}}{=} \sum_{l=1}^k \gamma_{nl}^{(l)} \in I\!\!R^n$. Because $\Sigma \succ 0$ it implies $\tilde{\gamma}_n^\top \Sigma \tilde{\gamma}_n > 0$ for all $\tilde{\gamma}_n \neq 0$. But even for $\gamma_{nk} \neq 0$, if its subvectors $\{\gamma_{nl}^{(l)}\}$ are linearly dependent, $\tilde{\gamma}_n$ can be zero. Thus there exists a nonzero vector γ such that $\gamma^\top (J_k \otimes \Sigma) \gamma = 0$. This means positive semidefiniteness.

The second assertion follows from the observation that Assumption (A4) due to the equality (4.13) also holds for the Kronecker product

$$(1 - \delta) J_k \otimes \Sigma \preceq J_k \otimes \Sigma_0 \preceq (1 + \delta) J_k \otimes \Sigma. \quad (4.14)$$

Therefore

$$(1 - \delta) \mathbf{D}_k (J_k \otimes \Sigma) \mathbf{D}_k^\top \preceq \mathbf{D}_k (J_k \otimes \Sigma_0) \mathbf{D}_k^\top \preceq (1 + \delta) \mathbf{D}_k (J_k \otimes \Sigma) \mathbf{D}_k^\top.$$

□

Lemma 4.9. Fix $x \in I\!\!R^d$. Suppose that the weights $\{w_{l,i}(x)\}$ satisfy

$$w_{l,i}(x) w_{m,i}(x) = w_{l,i}(x), \quad l \leq m. \quad (4.15)$$

Then under Assumptions (A1), (A2), (A5) the covariance matrix Σ_k defined by (3.5) is nonsingular with

$$\det \Sigma_k = \det \mathbf{B}_k^{-1} \prod_{l=2}^k \det(\mathbf{B}_{l-1}^{-1} - \mathbf{B}_l^{-1}) > 0, \quad k = 2, \dots, K. \quad (4.16)$$

Remark 4.2. The condition (4.15) holds for rectangular kernels with nested supports.

Proof. The condition (4.15) implies

$$\mathbf{W}_l \Sigma \mathbf{W}_m = \text{diag}(w_{l,1} w_{m,1} / \sigma_1^2, \dots, w_{l,n} w_{m,n} / \sigma_n^2) = \mathbf{W}_l$$

for any $l \leq m$. Thus the blocks of Σ_k simplify to

$$D_l \Sigma D_m^\top = \mathbf{B}_l^{-1} \Psi \mathbf{W}_l \Sigma \mathbf{W}_m \Psi^\top \mathbf{B}_m^{-1} = \mathbf{B}_l^{-1} \Psi \mathbf{W}_l \Psi^\top \mathbf{B}_m^{-1}$$

and Σ_k has a simple structure:

$$\Sigma_k = \begin{pmatrix} \mathbf{B}_1^{-1} & \mathbf{B}_2^{-1} & \mathbf{B}_3^{-1} & \dots & \mathbf{B}_k^{-1} \\ \mathbf{B}_2^{-1} & \mathbf{B}_2^{-1} & \mathbf{B}_3^{-1} & \dots & \mathbf{B}_k^{-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{B}_k^{-1} & \mathbf{B}_k^{-1} & \mathbf{B}_k^{-1} & \dots & \mathbf{B}_k^{-1} \end{pmatrix}.$$

Then the determinant of Σ_k coincides with the determinant of the following irreducible block triangular matrix:

$$\det \Sigma_k = \begin{vmatrix} \mathbf{B}_1^{-1} - \mathbf{B}_2^{-1} & \mathbf{B}_2^{-1} - \mathbf{B}_3^{-1} & \dots & \mathbf{B}_{k-1}^{-1} - \mathbf{B}_k^{-1} & \mathbf{B}_k^{-1} \\ \mathbf{0} & \mathbf{B}_2^{-1} - \mathbf{B}_3^{-1} & \dots & \mathbf{B}_{k-1}^{-1} - \mathbf{B}_k^{-1} & \mathbf{B}_k^{-1} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & \mathbf{B}_{k-1}^{-1} - \mathbf{B}_k^{-1} & \mathbf{B}_k^{-1} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{B}_k^{-1} \end{vmatrix}$$

implying

$$\det \Sigma_k = \det(\mathbf{B}_1^{-1} - \mathbf{B}_2^{-1}) \det(\mathbf{B}_2^{-1} - \mathbf{B}_3^{-1}) \cdots \det(\mathbf{B}_{k-1}^{-1} - \mathbf{B}_k^{-1}) \det \mathbf{B}_k^{-1}.$$

Clearly the matrix Σ_k is nonsingular if all the matrices $\mathbf{B}_{l-1}^{-1} - \mathbf{B}_l^{-1}$ are nonsingular. By (A1) and (A2) $\mathbf{B}_l \succ 0$ for any l . By (A5) there exists $u_0 > 1$ such that $\mathbf{B}_l \succeq u_0 \mathbf{B}_{l-1}$, therefore $\mathbf{B}_{l-1}^{-1} - \mathbf{B}_l^{-1} \succeq (1 - 1/u_0) \mathbf{B}_{l-1}^{-1} \succ \mathbf{B}_{l-1}^{-1} \succ 0$. \square

Lemma 4.10. *In the “nonparametric situation” the moment generation function (mgf) of the joint distribution of $\tilde{\boldsymbol{\theta}}_1, \dots, \tilde{\boldsymbol{\theta}}_K$ is*

$$\mathbb{E} \exp \left\{ \gamma^\top (\text{vec } \tilde{\boldsymbol{\Theta}}_K - \text{vec } \boldsymbol{\Theta}_K^*) \right\} = \exp \left\{ \frac{1}{2} \gamma^\top \Sigma_{K,0} \gamma \right\}. \quad (4.17)$$

Thus, provided that $\Sigma_{K,0} \succ 0$, it holds that $\text{vec } \tilde{\boldsymbol{\Theta}}_K \sim \mathcal{N}(\text{vec } \boldsymbol{\Theta}_K^*, \Sigma_{K,0})$.

Similarly, in the “parametric situation”, if $\Sigma_K \succ 0$, then the joint distribution of $\text{vec } \tilde{\boldsymbol{\Theta}}_K$ is $\mathcal{N}(\text{vec } \boldsymbol{\Theta}_K, \Sigma_K)$ with the mgf:

$$\mathbb{E} \exp \left\{ \gamma^\top (\text{vec } \tilde{\boldsymbol{\Theta}}_K - \text{vec } \boldsymbol{\Theta}_K) \right\} = \exp \left\{ \frac{1}{2} \gamma^\top \Sigma_K \gamma \right\}. \quad (4.18)$$

Proof. Let $\gamma \in \mathbb{R}^{pK}$ be written in a partitioned form $\gamma^\top = (\gamma_1^\top, \dots, \gamma_K^\top)$ with $\gamma_l \in \mathbb{R}^p$, $l = 1, \dots, K$. Then the mgf for the centered random vector $\text{vec } \tilde{\boldsymbol{\Theta}}_K - \text{vec } \boldsymbol{\Theta}_K^* \in \mathbb{R}^{pK}$, due to the decomposition (2.8) $\tilde{\boldsymbol{\theta}}_l = \boldsymbol{\theta}_l^* + D_l \Sigma_0^{1/2} \boldsymbol{\varepsilon}$ with $D_l = \mathbf{B}_l^{-1} \Psi \mathbf{W}_l$, can be represented as follows:

$$\begin{aligned} \mathbb{E} \exp \left\{ \gamma^\top (\text{vec } \tilde{\boldsymbol{\Theta}}_K - \text{vec } \boldsymbol{\Theta}_K^*) \right\} &= \mathbb{E} \exp \left\{ \sum_{l=1}^K \gamma_l^\top (\tilde{\boldsymbol{\theta}}_l - \boldsymbol{\theta}_l^*) \right\} \\ &= \mathbb{E} \exp \left\{ \sum_{l=1}^K \gamma_l^\top D_l \Sigma_0^{1/2} \boldsymbol{\varepsilon} \right\} = \mathbb{E} \exp \left\{ \left(\sum_{l=1}^K D_l^\top \gamma_l \right)^\top \Sigma_0^{1/2} \boldsymbol{\varepsilon} \right\}. \end{aligned}$$

A trivial observation that $\sum_{l=1}^K D_l^\top \gamma_l$ is a vector in \mathbb{R}^n and $\Sigma_0^{1/2} \boldsymbol{\varepsilon} \sim \mathcal{N}(0, \Sigma_0)$ by (1.1) implies by the definition of $\Sigma_{K,0}$ the first assertion of the lemma, because

$$\begin{aligned} \mathbb{E} \exp \left\{ \left(\sum_{l=1}^K D_l^\top \gamma_l \right)^\top \Sigma_0^{1/2} \boldsymbol{\varepsilon} \right\} &= \exp \left\{ \frac{1}{2} \left(\sum_{l=1}^K D_l^\top \gamma_l \right)^\top \Sigma_0 \left(\sum_{l=1}^K D_l^\top \gamma_l \right) \right\} \\ &= \exp \left\{ \frac{1}{2} (\mathbf{D}_K^\top \gamma)^\top (J_K \otimes \Sigma_0) \mathbf{D}_K^\top \gamma \right\} = \exp \left\{ \frac{1}{2} \gamma^\top \Sigma_{K,0} \gamma \right\}, \end{aligned}$$

here \mathbf{D}_K is defined by (3.23). \square

4.4 Proof of the propagation property

Lemma 4.11. *The Kullback-Leibler divergence between the distributions of $\text{vec } \tilde{\boldsymbol{\Theta}}_k$ under the alternative and under the null has the following form:*

$$\begin{aligned} 2\mathbb{KL}(\mathbb{P}_{\boldsymbol{f}, \Sigma_0}^k, \mathbb{P}_{\boldsymbol{\theta}, \Sigma}^k) &\stackrel{\text{def}}{=} 2\mathbb{E}_{\boldsymbol{f}, \Sigma_0} \log \left(\frac{d\mathbb{P}_{\boldsymbol{f}, \Sigma_0}^k}{d\mathbb{P}_{\boldsymbol{\theta}, \Sigma}^k} \right) \\ &= \Delta(k) + \log \left(\frac{\det \Sigma_k}{\det \Sigma_{k,0}} \right) + \text{tr}(\Sigma_k^{-1} \Sigma_{k,0}) - pk, \end{aligned} \quad (4.19)$$

where

$$b(k) \stackrel{\text{def}}{=} \text{vec } \boldsymbol{\Theta}_k^* - \text{vec } \boldsymbol{\Theta}_k \quad (4.20)$$

$$\Delta(k) \stackrel{\text{def}}{=} b(k)^\top \boldsymbol{\Sigma}_k^{-1} b(k). \quad (4.21)$$

Proof. Denote the Radon-Nikodym derivative by $Z_k \stackrel{\text{def}}{=} d\mathbb{P}_{f,\Sigma_0}^k / d\mathbb{P}_{\theta,\Sigma}^k$. Then

$$\begin{aligned} \log(Z_k(y)) &= \frac{1}{2} \log \left(\frac{\det \boldsymbol{\Sigma}_k}{\det \boldsymbol{\Sigma}_{k,0}} \right) - \frac{1}{2} \|\boldsymbol{\Sigma}_{k,0}^{-1/2}(y - \text{vec } \boldsymbol{\Theta}_k^*)\|^2 \\ &\quad + \frac{1}{2} \|\boldsymbol{\Sigma}_k^{-1/2}(y - \text{vec } \boldsymbol{\Theta}_k)\|^2 \end{aligned} \quad (4.22)$$

can be considered as a quadratic function of $\text{vec } \boldsymbol{\Theta}_k$. By the Taylor expansion at the point $\text{vec } \boldsymbol{\Theta}_k^*$ the last expression reads as follows

$$\begin{aligned} \log(Z_k(y)) &= \frac{1}{2} \log \left(\frac{\det \boldsymbol{\Sigma}_k}{\det \boldsymbol{\Sigma}_{k,0}} \right) - \frac{1}{2} \|\boldsymbol{\Sigma}_{k,0}^{-1/2}(y - \text{vec } \boldsymbol{\Theta}_k^*)\|^2 \\ &\quad + \frac{1}{2} \|\boldsymbol{\Sigma}_k^{-1/2}(y - \text{vec } \boldsymbol{\Theta}_k^*)\|^2 + b(k)^\top \boldsymbol{\Sigma}_k^{-1}(y - \text{vec } \boldsymbol{\Theta}_k^*) + \frac{1}{2} \Delta(k). \end{aligned}$$

Then the expression for the Kullback-Leibler divergence can be written in the following way:

$$\begin{aligned} \mathbb{KL}(\mathbb{P}_{f,\Sigma_0}^k, \mathbb{P}_{\theta,\Sigma}^k) &\stackrel{\text{def}}{=} \mathbb{E}_{f,\Sigma_0} \log(Z_k) \\ &= \frac{1}{2} \log \left(\frac{\det \boldsymbol{\Sigma}_k}{\det \boldsymbol{\Sigma}_{k,0}} \right) + \frac{1}{2} \Delta(k) + \frac{1}{2} \mathbb{E} \{ \|\boldsymbol{\Sigma}_k^{-1/2} \boldsymbol{\Sigma}_{k,0}^{1/2} \xi\|^2 - \|\xi\|^2 + 2b(k)^\top \boldsymbol{\Sigma}_k^{-1} \boldsymbol{\Sigma}_{k,0}^{1/2} \xi \}, \end{aligned}$$

where $\xi \sim \mathcal{N}(0, I_{pk})$. This implies

$$2\mathbb{KL}(\mathbb{P}_{f,\Sigma_0}^k, \mathbb{P}_{\theta,\Sigma}^k) = \Delta(k) + \log \left(\frac{\det \boldsymbol{\Sigma}_k}{\det \boldsymbol{\Sigma}_{k,0}} \right) + \text{tr}(\boldsymbol{\Sigma}_k^{-1} \boldsymbol{\Sigma}_{k,0}) - pk. \quad (4.23)$$

In the case of *homogeneous errors* with $\sigma_{0,i} = \sigma_0$ and $\sigma_i = \sigma, i = 1, \dots, n$ the calculations simplify a lot. Now

$$\boldsymbol{\Sigma}_k = \sigma^2 \mathbf{V}_k, \quad \boldsymbol{\Sigma}_{k,0} = \sigma_0^2 \mathbf{V}_k$$

with a $pk \times pk$ matrix \mathbf{V}_k defined as

$$\mathbf{V}_k = (\overline{D}_1 \oplus \dots \oplus \overline{D}_k)(J_k \otimes I_n)(\overline{D}_1 \oplus \dots \oplus \overline{D}_k)^\top,$$

where $\overline{D}_l = (\boldsymbol{\Psi} \mathcal{W}_l \boldsymbol{\Psi}^\top)^{-1} \boldsymbol{\Psi} \mathcal{W}_l$, $l = 1, \dots, k$ does not depend on σ . Then $\Delta(k) = \sigma^{-2} \Delta_1(k)$, with $\Delta_1(k) \stackrel{\text{def}}{=} b(k)^\top \mathbf{V}_k^{-1} b(k)$, $\det \boldsymbol{\Sigma}_k / \det \boldsymbol{\Sigma}_{k,0} = (\sigma^2 / \sigma_0^2)^{pk}$, and the expression for the Kullback-Leibler divergence reads as follows:

$$\begin{aligned} \mathbb{KL}(\mathbb{P}_{f,\Sigma_0}^k, \mathbb{P}_{\theta,\Sigma}^k) &= pk \log \left(\frac{\sigma}{\sigma_0} \right) + \frac{1}{2} \Delta(k) + \frac{pk}{2} \left(\frac{\sigma_0^2}{\sigma^2} - 1 \right) \\ &= pk \log \left(\frac{\sigma}{\sigma_0} \right) + \frac{1}{2\sigma^2} b(k)^\top \mathbf{V}_k^{-1} b(k) + \frac{pk}{2} \left(\frac{\sigma_0^2}{\sigma^2} - 1 \right), \end{aligned} \quad (4.24)$$

implying the same asymptotic behavior as in (3.13).

□

Proof. of Theorem 3.2 (Propagation property)

Notice that for any nonnegative measurable function $g = g(\tilde{\Theta}_k)$ the Cauchy-Schwarz inequality implies

$$\mathbb{E}_{f,\Sigma_0}[g] = \mathbb{E}_{\theta,\Sigma}[gZ_k] \leq (\mathbb{E}_{\theta,\Sigma}[g^2])^{1/2} (\mathbb{E}_{\theta,\Sigma}[Z_k^2])^{1/2} \quad (4.25)$$

with the Radon-Nikodym derivative $Z_k = dP_{f,\Sigma_0}^k / dP_{\theta,\Sigma}^k$. One gets the first assertion taking $g = |(\tilde{\theta}_k - \theta)^\top \mathbf{B}_k (\tilde{\theta}_k - \theta)|^{r/2}$, and applying “the parametric risk bound” with $\delta = 0$ from (4.7):

$$\begin{aligned} \mathbb{E}[g] &\leq (\mathbb{E}_{\theta,\Sigma}|(\tilde{\theta}_k - \theta)^\top \mathbf{B}_k (\tilde{\theta}_k - \theta)|^r)^{1/2} (\mathbb{E}_{\theta,\Sigma}[Z_k^2])^{1/2} \\ &= (\mathbb{E}_{\theta,\Sigma}|2 L(\mathbf{W}_k, \tilde{\theta}_k, \theta)|^r)^{1/2} (\mathbb{E}_{\theta,\Sigma}[Z_k^2])^{1/2} \\ &\leq (\mathbb{E}|\chi_p^2|^r)^{1/2} (\mathbb{E}_{\theta,\Sigma}[Z_k^2])^{1/2}. \end{aligned}$$

The second assertion is treated similarly by application of the pivotality property from Lemma 4.1 and the propagation conditions (2.17).

To calculate $\mathbb{E}_{\theta,\Sigma}[Z_k^2]$ let us consider $\log Z_k$ given by

$$\begin{aligned} \log(Z_k(y)) &= \frac{1}{2} \log \left(\frac{\det \Sigma_k}{\det \Sigma_{k,0}} \right) - \frac{1}{2} \|\Sigma_{k,0}^{-1/2}(y - \text{vec } \Theta_k^*)\|^2 \\ &\quad + \frac{1}{2} \|\Sigma_k^{-1/2}(y - \text{vec } \Theta_k)\|^2 \end{aligned}$$

as a function of $\text{vec } \Theta_k^*$. Application of the Taylor expansion at the point $\text{vec } \Theta_k$ yields

$$\begin{aligned} 2 \log Z_k &= \log \frac{\det \Sigma_k}{\det \Sigma_{k,0}} - \|\Sigma_{k,0}^{-1/2}(y - \text{vec } \Theta_k)\|^2 + \|\Sigma_k^{-1/2}(y - \text{vec } \Theta_k)\|^2 \\ &\quad + 2b(k)^\top \Sigma_{k,0}^{-1}(y - \text{vec } \Theta_k) - b(k)^\top \Sigma_{k,0}^{-1}b(k). \end{aligned}$$

With $\xi \sim \mathcal{N}(0, I_{pk})$ the second moment of the Radon-Nikodym derivative reads as follows

$$\begin{aligned} &\mathbb{E}_{\theta,\Sigma}[Z_k^2] \\ &= \frac{\det \Sigma_k}{\det \Sigma_{k,0}} \exp\{-b(k)^\top \Sigma_{k,0}^{-1}b(k)\} \mathbb{E} \exp\{-\|\Sigma_{k,0}^{-1/2}\Sigma_k^{1/2}\xi\|^2 + \|\xi\|^2 + 2b(k)^\top \Sigma_{k,0}^{-1}\Sigma_k^{1/2}\xi\} \\ &= \frac{\det \Sigma_k}{\det \Sigma_{k,0}} [\det(2\Sigma_k^{1/2}\Sigma_{k,0}^{-1}\Sigma_k^{1/2} - I_{pk})]^{-1/2} \\ &\times \exp\{2b(k)^\top \Sigma_{k,0}^{-1}\Sigma_k^{1/2}(2\Sigma_k^{1/2}\Sigma_{k,0}^{-1}\Sigma_k^{1/2} - I_{pk})^{-1}\Sigma_k^{1/2}\Sigma_{k,0}^{-1}b(k) - b(k)^\top \Sigma_{k,0}^{-1}b(k)\} \\ &= \frac{\det \Sigma_k}{\det \Sigma_{k,0}} \left[\prod_{j=1}^{pk} \{2\lambda_j(\Sigma_k^{1/2}\Sigma_{k,0}^{-1}\Sigma_k^{1/2}) - 1\} \right]^{-1/2} \quad (4.26) \\ &\times \exp\{b(k)^\top \Sigma_{k,0}^{-1/2} [2\Sigma_{k,0}^{-1/2}\Sigma_k^{1/2}(2\Sigma_k^{1/2}\Sigma_{k,0}^{-1}\Sigma_k^{1/2} - I_{pk})^{-1}\Sigma_k^{1/2}\Sigma_{k,0}^{-1/2} - I_{pk}] \Sigma_{k,0}^{-1/2}b(k)\}. \end{aligned}$$

To estimate the obtained expression in terms of the level of noise misspecification δ notice that the condition (3.8) implies

$$\left(\frac{1}{1+\delta} \right)^{pk} \leq \frac{\det \Sigma_k}{\det \Sigma_{k,0}} \leq \left(\frac{1}{1-\delta} \right)^{pk},$$

$$\begin{aligned} \left(\frac{1-\delta}{1+\delta}\right)^{\frac{pk}{2}} &\leq \left[\prod_{j=1}^{pk}\{2\lambda_j(\Sigma_k^{1/2}\Sigma_{k,0}^{-1}\Sigma_k^{1/2}) - 1\}\right]^{-1/2} \leq \left(\frac{1+\delta}{1-\delta}\right)^{\frac{pk}{2}}. \\ \frac{1-\delta}{1+\delta}I_{pk} &\preceq \left(2\Sigma_k^{1/2}\Sigma_{k,0}^{-1}\Sigma_k^{1/2} - I_{pk}\right)^{-1} \preceq \frac{1+\delta}{1-\delta}I_{pk}. \end{aligned}$$

Therefore the quantity in the exponent in (4.26) is bounded by:

$$\begin{aligned} &\left(2\frac{1-\delta}{(1+\delta)^2} - 1\right)b(k)^\top \Sigma_{k,0}^{-1}b(k) \\ &\leq b(k)^\top \Sigma_{k,0}^{-1/2} [2\Sigma_{k,0}^{-1/2}\Sigma_k^{1/2}(2\Sigma_k^{1/2}\Sigma_{k,0}^{-1}\Sigma_k^{1/2} - I_{pk})^{-1}\Sigma_k^{1/2}\Sigma_{k,0}^{-1/2} - I_{pk}] \Sigma_{k,0}^{-1/2}b(k) \\ &\leq \left(2\frac{1+\delta}{(1-\delta)^2} - 1\right)b(k)^\top \Sigma_{k,0}^{-1}b(k). \end{aligned}$$

Moreover,

$$\begin{aligned} \frac{\Delta(k)}{1+\delta} &= \frac{1}{1+\delta}b(k)^\top \Sigma_k^{-1}b(k) \\ &\leq b(k)^\top \Sigma_{k,0}^{-1}b(k) \\ &\leq \frac{1}{1-\delta}b(k)^\top \Sigma_k^{-1}b(k) = \frac{\Delta(k)}{1-\delta}. \end{aligned}$$

Finally,

$$\begin{aligned} &\left(\frac{1-\delta}{(1+\delta)^3}\right)^{\frac{pk}{2}} \exp\left\{\left(\frac{2(1-\delta)}{(1+\delta)^2} - 1\right)\frac{\Delta(k)}{1+\delta}\right\} \\ &\leq \mathbb{E}_{\boldsymbol{\theta},\Sigma}[Z_k^2] \leq \left(\frac{1+\delta}{(1-\delta)^3}\right)^{\frac{pk}{2}} \exp\left\{\left(\frac{2(1+\delta)}{(1-\delta)^2} - 1\right)\frac{\Delta(k)}{1-\delta}\right\}. \end{aligned} \quad (4.27)$$

In the case of homogeneous errors the expression for $\log Z_k$ reads as

$$\begin{aligned} \log Z_k &= pk \log\left(\frac{\sigma}{\sigma_0}\right) + \frac{1}{2}\left(\frac{1}{\sigma^2} - \frac{1}{\sigma_0^2}\right)\|\mathbf{V}_k^{-1/2}(y - \text{vec } \boldsymbol{\Theta}_k)\|^2 \\ &+ \frac{1}{\sigma_0^2}b(k)^\top \mathbf{V}_k^{-1}(y - \text{vec } \boldsymbol{\Theta}_k) - \frac{1}{2\sigma_0^2}b(k)^\top \mathbf{V}_k^{-1}b(k), \end{aligned}$$

implying

$$\mathbb{E}_{\boldsymbol{\theta},\sigma}[Z_k^2] = \left(\frac{\sigma^2}{\sigma_0^2}\right)^{pk} \left(\frac{\sigma_0^2}{2\sigma^2 - \sigma_0^2}\right)^{\frac{pk}{2}} \exp\left\{\frac{b(k)^\top \mathbf{V}_k^{-1}b(k)}{2\sigma^2 - \sigma_0^2}\right\}.$$

By Assumption (A4)

$$\begin{aligned} &\left(\frac{1-\delta}{(1+\delta)^3}\right)^{\frac{pk}{2}} \exp\left\{\frac{\Delta_1(k)}{\sigma^2(1+\delta)}\right\} \\ &\leq \mathbb{E}_{\boldsymbol{\theta},\sigma}[Z_k^2] \leq \left(\frac{1+\delta}{(1-\delta)^3}\right)^{\frac{pk}{2}} \exp\left\{\frac{\Delta_1(k)}{\sigma^2(1-\delta)}\right\}, \end{aligned} \quad (4.28)$$

where p is the dimension of the parameter set and k is the degree of the localization. \square

References

- [1] Akaike, H. (1973). Information theory and an extension of the maximum likelihood principle. *Second International Symposium on Information Theory (Tsahkadsor, 1971), Akadmiai Kiad, Budapest* 267–281.
- [2] Baraud, Y., Giraud,C. and Huet,S. (2009) Gaussian model selection with an unknown variance. *Ann. Statist.* **37:2** 630–672.
- [3] Birgé, L. and Massart, P. (2001) Gaussian model selection. *Journal of the European Mathematical Society* **3:3** 203–268.
- [4] Brua, J.-Y. (2009). Asymptotic efficient estimators for non-parametric heteroscedastic model. *Statistical Methodology* **6:1** 47–60.
- [5] Donoho, D. L. and Johnstone, I. M. (1994). Ideal spatial adaptation by wavelet shrinkage. *Biometrika* **81** 425–455
- [6] Fan, J. and Gijbels, I. (1996). *Local Polynomial Modelling and Its Applications*. Monographs on Statistics and Applied Probability, 66. Chapman and Hall, London.
- [7] Fan, J., Zhang, C. and Zhang, J. (2001). Generalized likelihood ratio statistics and Wilks phenomenon. *Ann. Statist.* **29:1** 153–193.
- [8] Galtchouk, L. and Pergamenshchikov, S. (2009). Adaptive asymptotically efficient estimation in heteroscedastic nonparametric regression *Journal of the Korean Statistical Society* **38:4** 305–322.
- [9] Galtchouk, L. and Pergamenshchikov, S. (2010). Adaptive asymptotically efficient estimation in heteroscedastic nonparametric regression. arXiv:1002.1537v1
- [10] Galtchouk, L. and Pergamenshchikov, S. (2010). Sharp non-asymptotic oracle inequalities for nonparametric heteroscedastic regression models. arXiv:1002.1538v1
- [11] Goldenshluger, A. and Nemirovski, A. (1994). On spatial adaptive estimation of nonparametric regression. *Research report, Technion-Israel Inst. Technology, Haifa, Israel*.
- [12] Efroimovich, S. and Pinsker, M.(1996). Sharp-optimal and adaptive estimation for heteroscedastic nonparametric regression. *Statistica Sinica* **6** 925–942.
- [13] Efroimovich, S. (2007). Sequential design and estimation in heteroscedastic nonparametric regression. *Sequential Analysis* **26** 3–25.
- [14] Katkovnik, V. Ja. (1979). Linear and nonlinear methods of nonparametric regression analysis. (Russian) *Soviet Automat. Control* **5** 35–46, 93.
- [15] Katkovnik, V. Ja. (1983). Convergence of linear and nonlinear nonparametric estimates of “kernel” type. *Automat. Remote Control* **44:4** 495–506; translated from *Avtomat. i Telemekh.* 1983 **4** 108–120 (Russian).
- [16] Katkovnik, V. Ja. (1985). *Nonparametric Identification and Data Smoothing: Local Approximation Approach*. Nauka, Moscow (Russian).

- [17] Katkovnik, V., Egiazarian, K. and Astola, J. (2006). *Local Approximation Techniques in Signal and Image Processing*. Bellingham, WA: SPIE Press.
- [18] Katkovnik, V. and Spokoiny, V. (2008). Spatially adaptive estimation via fitted local likelihood techniques. *IEEE Trans. Signal Process.*, **56**, No.3, 873–886.
- [19] Kullback, S. and Leibler, R. A. (1951). On information and sufficiency. *Ann. Math. Statistics* **22** 79–86.
- [20] Lepskii, O. V. (1990). A problem of adaptive estimation in Gaussian white noise. (Russian) *Teor. Veroyatnost. i Primenen.* **35:3** 459–470; translation in *Theory Probab. Appl.* **35:3** 454–466.
- [21] Lepskii, O. V. (1992). Asymptotic minimax adaptive estimation. II. Schemes without optimal adaptation. Adaptive estimates. (Russian) *Teor. Veroyatnost. i Primenen.* **37:3** 468–481; translation in *Theory Probab. Appl.* **37:3** 433–448.
- [22] Lepski, O. V., Mammen, E. and Spokoiny, V.G. (1997). Optimal spatial adaptation to inhomogeneous smoothness: an approach based on kernel estimates with variable bandwidth selectors. *Ann. Stat.* **25:3** 929–947.
- [23] Lepski, O. V. and Spokoiny, V.G. (1997). Optimal pointwise adaptive methods in nonparametric estimation. *Ann. Stat.* **25:6** 2512–2546.
- [24] Loader, C. (1999). *Local Regression and Likelihood*. Statistics and Computing. Springer-Verlag, New York.
- [25] Massart, P. (2003) *Concentration Inequalities and Model Selection* (2007). Ecole d’été de Probabilités de Saint-Flour . Lecture Notes in Mathematics 1896, Springer Berlin/Heidelberg.
- [26] Spokoiny, V. (2002). Variance estimation for high-dimensional regression models. *J. Multivariate Anal.* **82** 111–133.
- [27] Spokoiny, V. and Vial, C. (2009). Parameter tuning in pointwise adaptation using a propagation approach. *Ann. Statist.* **37:5B** 2783–2807.
- [28] Tsybakov, A. B. (2009). *Introduction to Nonparametric Estimation*. Springer Series in Statistics. Springer-Verlag, New York.
- [29] White, H. (1982). Maximum likelihood estimation of misspecified models. *Econometrica* **50:1** 1–25.